

ТАРТУСКИЙ
ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ



ТРУДЫ

ВЫЧИСЛИТЕЛЬНОГО ЦЕНТРА

44

ТАРТУ
1980

ТАРТУСКИЙ ГОСУДАРСТВЕННЫЙ
УНИВЕРСИТЕТ

МЕТОДЫ МНОГОМЕРНОГО
СТАТИСТИЧЕСКОГО АНАЛИЗА

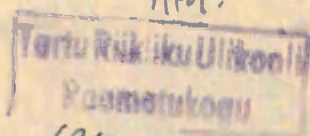
ТРУДЫ
ВЫЧИСЛИТЕЛЬНОГО
ЦЕНТРА

ВЫПУСК 44

ТАРТУ 1980

Утверждено на заседании совета математического
факультета ТГУ 15 ноября 1979 года.

Ам.



6344 6296

KUSTUTATUD

КЛАССИФИКАЦИЯ И КАНОНИЧЕСКИЕ МОДЕЛИ СБАЛАНСИРОВАННЫХ ПОЛНЫХ ФАКТОРНЫХ КОМПЛЕКСОВ С ПРОИЗВОЛЬНЫМ ЧИСЛОМ ФАКТОРОВ

Т. Мелс

1. Введение

В настоящей статье описаны в матричной символике некоторые модели сбалансированных дисперсионных анализов (ДА). Поскольку такое исследование в хорошо проработанной классической области может показаться странным, необходимо некоторое обоснование тематики статьи.

Начнем обращением внимания на то, что при большом числе публикаций по линейным статистическим методам, только небольшое число из них занимается систематикой и описанием целых классов ДА. В руководствах (например [3] — [7]) подробно описаны многие частные схемы ДА, однако о многих других возможных схемах иногда даже не упоминается. Поэтому нередко трудно найти в готовом виде подходящую схему ДА для обработки имеющихся данных. Такая ситуация хорошо известна каждому консультанту по математической статистике. Нам кажется, что хорошая систематика ДА может заметно повышать эффективность применения ДА и избегать ошибок, связанных с использованием неадекватных схем.

Есть у проблемы и чисто вычислительный аспект. В много-

факторных моделях ДА численное решение нормальных уравнений может оказаться трудоемкой задачей, особенно там, где нет достаточно мощных ЭВМ. Если, однако, известно аналитическое решение нормальных уравнений, как например в сбалансированных случаях ДА, то объем необходимых вычислений может быть заметно сокращен. Поэтому желательно иметь в библиотеке СП рядом с общими модулями ДА также и специальные модули для разных классов сбалансированных ДА. Здесь, конечно, необходимы специализированные средства описания схем ДА, что требует соответствующих теоретических исследований, в первую очередь по систематике схем ДА.

В настоящей статье сделана попытка дать общее определение и описание т.н. полных факторных комплексов (ПФК). Для сбалансированных ПФК предлагается также теоретический вывод таблицы ДА. Развита соответствующая матричная техника.

В класс ПФК входят хорошо известные "факторные планы", планы с группировкой и их разные гибриды. Неполные факторные комплексы (латинские квадраты и др.) в статье не рассматриваются, хотя и они могут быть довольно элегантно включены в общий формализм факторных комплексов (см. 2.7).

2. Полные факторные комплексы

В этом параграфе определяем полные факторные комплексы и некоторые смежные понятия. Даем также таблицу всевозможных типов ПФК, если число факторов не превышает пяти.

2.1. Любой фактор будем представлять множеством уровней. Например, если фактор A имеет уровни a_1, \dots, a_s , то $A =$

$= \{a_1, \dots, a_n\}$. Каждому индивиду (=наблюдению, =измерению) $\omega \in \Omega$ (Ω - совокупность индивидов) соответствует некоторый элемент $\langle a, b, \dots, d, e \rangle$ во множестве $M = A \times B \times \dots \times D \times E$, где A, \dots, E - факторы, а a, \dots, e - их уровни, потенциально оказывающие влияние на индивида ω . Совокупность Ω всех индивидов определяет таким путем в M некоторое подмножество Z , которое называем структурой (совокупности индивидов).

2.2. Будем говорить, что две уровни комбинируются (в выбранной совокупности Ω), если они являются координатами некоторой точки структуры S . Если уровни комбинируются, то в Ω найдется индивид, на котором эти уровни реализуются одновременно.

Если совокупность индивидов не пуста, то у всех факторов найдутся уровни, которые комбинируются. Однако мы будем всегда предполагать, что у каждого фактора все уровни комбинируются с некоторыми уровнями других факторов: уровни, которые не комбинируются, можно просто отбросить. В силу сделанного соглашения проекция структуры S на любое координатное множество совпадает с этим множеством.

2.3. Рассматриваем два специальные типа связи между двумя факторами A и B . Будем говорить, что фактор A подчиняется фактору B (символически $A > B$ или $B < A$), если каждый уровень фактора A комбинируется (в Ω) только с одним уровнем фактора B . Если $A < B$ и $B < A$, то факторы A и B являются копиями друг друга и мы называем их эквивалентными. Все эквивалентные факторы могут быть отождествлены без всякого ущерба для теории. Поэтому далее будем всегда предполагать,

что любой фактор эквивалентен только самому себе. В силу этого допущения множество факторов вместе с бинарным отношением \prec превратится в частично упорядоченное множество.

Обозначим $\varphi_1, \dots, \varphi_\tau$ всевозможные частичные упорядочения фиксированного множества факторов. Каждое упорядочение φ_1 удобно интерпретировать как предикат: $\varphi_1(s)$ означает, что структура s порождает на множестве факторов упорядочение φ_1 . Предикаты φ_1 называем предикатами подчинения.

2.4. Другой специальный тип связи между двумя факторами — это перекрестная связь. Будем говорить, что факторы A и B связаны перекрестно (символически $A \times B$ или $B \times A$), если каждый уровень фактора A комбинируется (в Ω) с каждым уровнем фактора B . Предикаты всевозможных перекрестных связей на некотором фиксированном множестве факторов обозначим ψ_1, \dots, ψ_τ . Здесь $\psi_1(s)$ означает, что в структуре s между факторами имеются перекрестные связи ψ_1 .

2.5. При фиксированном множестве факторов множество всевозможных структур частично упорядочено теоретико-множественным включением \subset . Включение $s_1 \subset s_2$ означает, что все комбинации уровней, составляющие структуру s_1 , имеются также в структуре s_2 .

Пусть φ — некоторый предикат подчинения. Рассматриваем максимальные элементы подмножества $\mathcal{S}_\varphi = \{s : \varphi(s)\}$ в частично упорядоченном множестве всех структур (на M). Каждый максимальный элемент (структура) множества \mathcal{S}_φ устанавливает между факторами связи, где имеются всевозможные комбинации уровней, не противоречащие подчинениям φ . Мы будем называть

эти максимальные структуры полными факторными комплексами (ПФК).

2.6. Приведем пример трехфакторного ПФК. Пусть $A = \{a_1, a_2\}$, $B = \{b_1, b_2, b_3\}$, $C = \{c_1, c_2, c_3, c_4\}$, а φ состоит из подчинений $A < B$ и $A < C$. Тогда структура

$$S = \{\langle a_1, b_1, c_1 \rangle, \langle a_1, b_1, c_2 \rangle, \langle a_1, b_2, c_1 \rangle, \\ \langle a_1, b_2, c_2 \rangle, \langle a_2, b_3, c_3 \rangle, \langle a_2, b_3, c_4 \rangle\}$$


является ПФК. Действительно, структура S порождает как подчинение $A < B$ (поскольку b_1 и b_2 соответствует только a_1 , а b_3 соответствует только a_2), так и подчинение $A < C$ ($c_1, c_2 \rightarrow a_1$, $c_3, c_4 \rightarrow a_2$). Но никакая более мощная структура на $A \times B \times C$ этими подчинениями уже не обладает.


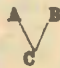

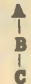
2.7. Аналогично ПФК могут быть определены некоторые дуальные им структуры как минимальные (в смысле отношения \subset) элементы множества $\mathcal{F}_\varphi = \{S : \varphi(S)\}$, где φ — предикат перекрестных связей. Например, латинский квадрат может быть определен как минимальный элемент среди трехфакторных структур, где каждые два фактора связаны перекрестно. Комбинируя, далее, переходы к минимальным и максимальным элементам, можно получить огромное количество интересных факторных структур. Эти вопросы не являются, однако, предметом настоящей статьи.


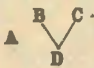
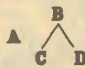
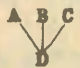


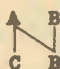

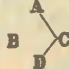
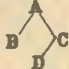
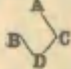




2.8. Переходим к понятию типа ПФК. Ясно, что в общих чертах тип ПФК определяется предикатом подчинения или (что равносильно) соответствующим частичным упорядочением это—


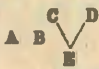
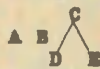
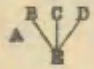
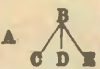
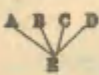

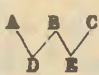
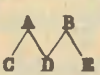
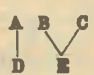

Т а б л и ц а I

Структурные графы полных факторных комплексов

Двухфакторные комплексы					
A B					

Трехфакторные комплексы					
A B C					

Четырехфакторные комплексы					
A B C D					
					
					

Пятифакторные комплексы					
A B C D E					
					

Пятифакторные комплексы (продолж.)

го ПФК. Но разные упорядочения φ_1 на множестве факторов могут быть изоморфны в том смысле, что после подходящего переименования факторов они совпадают. Изоморфны, например, упорядочения $\varphi_1 = \{A < B\}$ и $\varphi_2 = \{B < A\}$. Поэтому типом ПФК естественно называть не само упорядочение φ этого ПФК, а класс эквивалентности упорядочения φ по изоморфности.

Непосредственным перебором вариантов можно убедиться, что в случае двух факторов A и B имеется два типа ПФК — $\{ \}$ и $\{A < B\}$. В случае трех факторов имеется уже пять типов ПФК: $\{ \}$, $\{A < B\}$, $\{A < B, A < C\}$, $\{A < B, C < B\}$ и $\{A < B < C\}$. Для четырех- и пятифакторных ПФК различных типов соответственно 16 и 63.

Конкретные примеры на все типы двух-, трех- и четырехфакторных ПФК приведены в [2]. Эти примеры показывают, что всевозможные типы ПФК могут вполне естественным образом встречаться в исследованиях по биологии, технике и медицине.

2.9. Различные типы ПФК наглядно и удобно представлять ориентированными графами: вершины графа изображают факторы, а каждому подчинению $A < B$ соответствует ребро, направленное из A в B . Если начертить эти графы таким образом, чтобы в них горизонтальных ребр нет, а все ребра считать направленными сверху вниз (тогда на ребрах стрелки указывать не надо), то получим "графы", которые называем структурными графами. Структурные графы всевозможных типов ПФК с числом факторов до пяти изображены в таблице 1.

2.10. Выделяем теперь ПФК, обладающие определенной симметрией. Пусть $A = \{a_1, \dots, a_s\}$ — фактор в ПФК. Обозначим λ

класс всех факторов, которым A подчиняется. Тогда любой уровень $a_1 \in A$ определяет однозначно некоторый кортеж $f(a_1)$ уровней факторов класса A . Введем в A отношение эквивалентности ("эквивалентность подчинения") \sim , считая $a_1 \sim a_j$, если $f(a_1) = f(a_j)$. Если, например, фактор A не подчиняется другим факторам (т.е. $A = \emptyset$), то $f(a_1)$ — пустой кортеж и при всех $a_1, a_j \in A$ имеем $a_1 \sim a_j$.

Разбиение A на классы эквивалентности по \sim называем естественным разбиением фактора. ПФК называем симметричным, если классы естественного разбиения любого фактора всегда имеют одинаковую мощность (зависящую от фактора). Симметричны, например, все ПФК без подчинений.

Симметричный ПФК называем сбалансированным, если все имеющиеся в нем комбинации факторных уровней реализуются на одинаковом числе (≥ 1) индивидов. Основным объектом дальнейшего изучения в статье и являются сбалансированные ПФК.

2.11. В случае ПФК удобно пользоваться т.н. сокращенной нумерацией уровней факторов, которая широко используется в разных вариантах ДА. Техника применения этой методики заключается в следующем. Каждому фактору в ПФК сопоставляется (латинский) индекс, называемый коренным индексом фактора. Внутри каждого естественного класса эквивалентности уровни фактора нумеруются коренным индексом фактора. Класс эквивалентности в свою очередь характеризуется однозначно значениями коренных индексов всех факторов, которым данный фактор подчиняется. Эти индексы указываются в скобках за коренным индексом фактора и называются дополнительными индексами.

Если ПФК симметричен, то у всякого фактора все классы эквивалентности имеют равную мощность, которую обозначаем n_A у фактора А, n_B у фактора В и т.д. и называем сокращенной размерностью фактора. Коренный индекс фактора А в симметричном ПФК пробегает значения от 1 до n_A .

В сбалансированном ПФК постоянным является также и число индивидов при каждой комбинации факторных уровней; это число (называемое повторностью измерения) обозначим n , а коренный индекс отдельного индивида (наблюдения) обозначим v ($v=1, \dots, n$). Общее число индивидов в совокупности Ω , структура которой является сбалансированным ПФК, равно произведению всех сокращенных размерностей факторов ПФК на повторность n .

3. Линейные модели сбалансированных ПФК

3.1. Переходим к конкретизации линейной модели

$$Y = K\psi + \varepsilon \quad (1)$$

для некоторого сбалансированного k -факторного ПФК. Как обычно, Y является здесь $(N \times q)$ -матрицей наблюдений (N — число индивидов в Ω , q — размерность признак-вектора), а матрицы K и ψ называются соответственно плановой матрицей и параметрической матрицей. Остаток ε считаем случайной центрированной матрицей. .

Конкретизация модели (1) состоит в установлении порядка, в котором стоят численные значения признаков в матрице Y , и в выборе параметрической матрицы ψ и соответствующей плановой матрицы K . В случае сбалансированных ПФК некоторый канонический

нический метод этих конкретизаций подсказывается структурой ПФК. Несколько следующих пунктов посвящены описанию этого метода.

3.2. Результаты измерения q исследуемых признаков у индивида образуют строку в Y . Поэтому строки матрицы наблюдений можно упорядочить как индивиды в X . Каноническим является следующий способ. Сначала упорядочиваем некоторым образом (обычно в алфавитном порядке) множество всех коренных индексов, включая ν , который занимает последнее место. Выбранный порядок индексов фиксируем и будем называть каноническим. Набор $k+1$ индексов, взятых в каноническом порядке, определяет однозначно индивид, а значит, и строку в Y . Если теперь упорядочить значения индексных наборов лексикографически, то и получим искомый канонический вид матрицы наблюдений.

3.3. Чтобы конкретизировать правую часть (1), будем выделить в K некоторые блоки столбцов, $K = (K_1 \vdots \dots \vdots K_r)$, а в ψ выделим соответствующие им блоки строк ψ_1 . Модель (1) может быть тогда представлен в виде

$$Y = \sum_{i=1}^r K_i \psi_i + \epsilon. \quad (2)$$

Рассматриваем естественный (канонический) выбор блоков K_1 и ψ_1 в случае сбалансированного ПФК. Все блоки в параметрической матрице ψ разделяются на три класса. Первый класс состоит только из одного блока ψ_1 — вектора общих средних всех q признаков:

$$\psi_1 = (\mu^{(1)}, \dots, \mu^{(q)}) .$$

Второй класс образуют блоки т.н. главных эффектов уровней факторов: главные эффекты каждого фактора образуют отдельный блок. Каждому уровню фактора соответствует q главных эффектов – на каждый признак один эффект. Эти эффекты образуют строку блока.

Главные эффекты обозначим одинаково с уровнями фактора, используя сокращенную нумерацию. При необходимости указываем верхним индексом еще номер признака (т.е. номер столбца в блоке). Например, $a_{i(k)}^{(t)}$ есть главный эффект факторного уровня $a_{i(k)}$ на t -тый признак. Строки в блоке упорядочиваем лексикографически по набору нижних индексов. Например, блок главных эффектов вида $a_{i(k)}^{(t)}$ может выглядеть как

$$\psi_2 = \begin{pmatrix} a_{1(1)}^{(1)} & a_{1(1)}^{(2)} \\ a_{1(2)}^{(1)} & a_{1(2)}^{(2)} \\ a_{2(1)}^{(1)} & a_{2(1)}^{(2)} \\ a_{2(2)}^{(1)} & a_{2(2)}^{(2)} \end{pmatrix}.$$

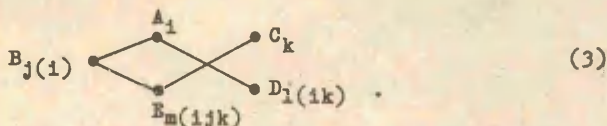
Третий класс образуют блоки из эффектов взаимодействий уровней разных факторов. Каждому набору факторов, где ни один фактор не подчиняется никакому другому фактору набора, соответствует блок взаимодействий. Блоки взаимодействий построены аналогично блокам главных эффектов. Например, эффект взаимодействия уровней $a_{i(k)}$ и $b_{j(k_1)}$ на t -тый признак обозначим $ab_{ij(k_1)}^{(t)}$ и этот эффект находится в t -том столбце и $ij(k_1)$ -ой строке блока взаимодействий факторов A и B .

3.4. Описываем блоки K_1 , соответствующие каноническому разложению параметрической матрицы. Введем для этого обозначения

- 1_p — единичная матрица порядка p ,
- $\dot{1}_p$ — столбец из p единиц,
- $\check{1}_p$ — специальная матрица, $\check{1}_p^T = (1_{p-1} : -\dot{1}_{p-1})$.

Каждый блок K_1 имеет вид тензорного произведения $k+1$ матриц 1_p , $\dot{1}_p$ или $\check{1}_p$ (k — число факторов в ПФК). Состав сомножителей в K_1 определяется нижними индексами элементов блока ψ_1 , а порядок сомножителей — установленным каноническим порядком индексов. Каждому индексу, который не встречается у элементов блока ψ_1 , соответствует в K_1 (на позиции, фиксированной для этого индекса) сомножитель $\dot{1}_p$. Коренным индексам ψ_1 соответствуют в K_1 сомножители $\check{1}_p$, а дополнительным индексам — сомножители 1_p . Размерность p каждого сомножителя в K_1 равно сокращенной размерности фактора, соответствующего этому индексу.

3.5. Рассмотрим пример модели пятифакторного сбалансированного ПФК. Пусть факторы А–Е связаны в ПФК согласно структурному графу



На графе указаны также коренные и дополнительные индексы факторов. Фиксируем следующий порядок индексов: i, j, k, l, m, v .

В данном случае параметрическая матрица разлагается на 10 блоков ψ_1, \dots, ψ_{10} с элементами соответственно $\mu^{(t)}$, $a_1^{(t)}$,

$b_{j(1)}^{(t)}, c_k^{(t)}, d_{l(1k)}^{(t)}, e_{m(1jk)}^{(t)}, a_{1k}^{(t)}, b_{jk(1)}^{(t)}, d_{jl(1k)}^{(t)}$ и $ed_{lm(1jk)}^{(t)}$. Взаимодействия трех или большего числа факторов в модели данного ПФК отсутствуют.

Канонический вид плановой матрицы получается следующий:

$$K = (\begin{matrix} 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{matrix}), \quad (4)$$

где сомножители в каждом блоке имеют соответственно n_A, n_B, n_C, n_D, n_E и n строк.

Если по найденным \downarrow и K вычислить произвольный элемент правой части (1), то получим для него равенство

$$y_{ijklm}^{(t)} = \mu^{(t)} + a_1^{(t)} + b_{j(1)}^{(t)} + c_k^{(t)} + d_{l(1k)}^{(t)} + e_{m(1jk)}^{(t)} + a_{1k}^{(t)} + b_{jk(1)}^{(t)} + d_{jl(1k)}^{(t)} + ed_{lm(1jk)}^{(t)} + \varepsilon_{ijklm}^{(t)}, \quad (5)$$

где $i=1, \dots, n_A; j=1, \dots, n_B; \dots; v=1, \dots, n$. Равенство (5) дает представление модели (2) в компонентах.

3.6. С помощью канонической плановой матрицы нетрудно найти проектор Q , порождающий т.н. остаточную сумму квадратов $Y^T Q Y$. Известно [5], что

$$Q = I - K(K^T K)^{-1} K^T, \quad (6)$$

где $^{-1}$ обозначает обобщенную обратную матрицу. Вычислить $K^T K$ помогает следующая лемма.

Лемма 1. Если K_1 и K_2 - разные канонические блоки в плановой матрице сбалансированного ПФК, то $K_1^T K_2 = 0$.

Доказательство. Обозначим $K_1 = \alpha_1 \otimes \dots \otimes \alpha_{k+1}$ и $K_j = \beta_1 \otimes \dots \otimes \beta_{k+1}$, где $\alpha_1, \dots, \beta_{k+1}$ равны 1, $\dot{1}$ или $\bar{1}$ (размерность матрицы зависит от порядкового номера сомножителя). Тогда

$$K_1^T K_j = \alpha_1^T \beta_1 \otimes \dots \otimes \alpha_{k+1}^T \beta_{k+1}. \quad (7)$$

Выражение (7) равно нулю, коль скоро в нем некоторый сомножитель $\alpha_1^T \beta_1$ имеет вид $\bar{1}^T \dot{1}$, так как

$$\bar{1}^T \dot{1} = 0. \quad (8)$$

Обозначим множества коренных и дополнительных индексов блока ψ_1 соответственно через G и J_G , а у блока ψ_j соответственно через F и J_F . Пусть, кроме того, $F = F_1 \cup H$ и $G = G_1 \cup H$, где $F_1 \cap G = G_1 \cap F = \emptyset$. Тогда, в силу (8), $K_1^T K_j$ может не равняться нулю лишь при условии, что каждому F_1 -сомножителю ($= \dot{1}$) в K_1 соответствует в K_j сомножитель 1, т.е. J_G -сомножитель. Другими словами, необходимым условием для $K_1^T K_j \neq 0$ является $F_1 \subset J_G$.

Включение $F_1 \subset J_G$ означает, что каждому F_1 -фактору подчиняется некоторый G -фактор. Поскольку никакой F -фактор не подчиняется никакому отличному от него F_1 -фактору ($F_1 \subset F$, а F — множество коренных индексов эффекта!), то G -факторы, которые подчиняются F_1 -факторам, не входят в F . Значит, они входят в $G - F = G_1$. Увидим, что если $K_1^T K_j \neq 0$, то каждому F_1 -фактору подчиняется некоторый G_1 -фактор. Но наше рассуждение вполне симметрично относительно G_1 и F_1 . Поэтому верно также, что при $K_1^T K_j \neq 0$ каждому G_1 -фактору подчиняется некоторый F_1 -фактор. Стало быть $F_1 \cap G_1 \neq \emptyset$, что, однако, невозможно. Лемма доказана.

На основе доказанной леммы, K^TK имеет блочно-диагональный вид, где на диагонали находятся блоки $K_1^TK_1$, являющиеся, очевидно, тензорными произведениями матриц i^Ti , $\check{i}^T\check{i}$ или 1. Эти матрицы регулярные, так что $(K_1^TK_1)^- = (K_1^TK_1)^{-1}$. Поэтому $(K^TK)^- = (K^TK)^{-1}$ и $(K^TK)^{-1}$ также имеет блочно-диагональный вид, где диагональные блоки являются тензорными произведениями матриц $(i^Ti)^{-1}$, $(\check{i}^T\check{i})^{-1}$ или 1. Теперь легко найти, что

$$Q = 1 - \sum_{t=1}^T K_t (K_t^TK_t)^{-1} K_t^T. \quad (9)$$

3.7. Если, согласно некоторой нулевой гипотезе $H_{0,t}$, t -тый эффект отсутствует (= нулевой), то блоки ψ_t и K_t в (2) следует пропустить. Проектор Q принимает при $H_{0,t}$ вид

$$Q_t = Q + K_t (K_t^TK_t)^{-1} K_t^T,$$

а остаточная сумма квадратов увеличивается на

$$y^T(Q_t - Q)y = y^TK_t(K_t^TK_t)^{-1}K_t^Ty. \quad (10)$$

Выясним более подробно структуру проектора $K_t(K_t^TK_t)^{-1}K_t^T$. Мы уже установили, что эта матрица является тензорным произведением $k+1$ сомножителей трех возможных типов, соответствующих трем возможным ролям индексов в элементах ψ_t . Рассмотрим некоторый индекс и пусть p - сокращенная размерность его фактора. Если t -тый блок ψ_t не содержит данного индекса, то индексу соответствует в тензорном произведении сомножитель

$$i_p(i_p^Ti_p)^{-1}i_p^T = p^{-1}U_p \equiv \bar{U}_p, \quad (11)$$

где U_p - $(p \times p)$ -матрица из единиц. Далее, если данный индекс является коренным индексом эффекта, то ему соответствует сомножитель вида

$$\check{1}_p (\check{1}_p^T \check{1}_p)^{-1} \check{1}_p^T = \check{1}_p (1_{p-1} - \frac{1}{p} U_{p-1}) \check{1}_p^T = 1_p - \bar{U}_p, \quad (12)$$

а любому добавочному индексу соответствует сомножитель 1_p . Например, для блока $\psi_t = \psi_9 = (bd_{jl(ik)})$ в модели (5) проектор в (10) имеет строение

$$Q_t - Q = 1_{n_A} \otimes (1_{n_B} - \bar{U}_{n_B}) \otimes 1_{n_C} \otimes (1_{n_D} - \bar{U}_{n_D}) \otimes \bar{U}_{n_E} \otimes \bar{U}_n. \quad (13)$$

3.8. Возвращаемся к вычислению (10). Основная трудность здесь заключается в несовместимости тензорного характера проектора Q и нетензорной природы матрицы U . Чтобы обойти это затруднение, преобразуем матрицу наблюдений к тензорному виду, используя формальные факторные операторы \hat{A}, \dots, \hat{E} и $\hat{\eta}$, а также формальный индивид σ . Будем предполагать, что факторные операторы и индивид σ являются элементами некоторых векторных пространств, которые имеют соответственно размерности n_A, \dots, n_E, n и q (q - размерность исследуемого в ДА вектора признаков).

Компоненты (нечисловые "координаты") операторов $\hat{A}, \dots, \hat{\eta}$ в некотором ортогональном базисе обозначим $\hat{A}_1, \dots, \hat{A}_{n_A}; \dots; \hat{\eta}_1, \dots, \hat{\eta}_n$. Сочетания $\hat{A}_1 \hat{B}_j \dots \hat{\eta}_v$ называем операторами измерения. Мы определяем их правилом

$$(\hat{A}_1 \dots \hat{E}_m \hat{\eta}_v)(\sigma_1, \dots, \sigma_q) = (y_{1 \dots m v}^{(1)}, \dots, y_{1 \dots m v}^{(q)}), \quad (14)$$

где $(\sigma_1, \dots, \sigma_q) = \sigma$, а $\sigma_1, \dots, \sigma_q$ удобно называть обобщенными

признаками. Оператор измерения, действуя на обобщенный признак, порождает число. Эту операцию естественно интерпретировать как генерирование значения признака зависящим от оператора способом.

3.9. Алгебраические и матричные операции над результатами измерения сводятся теперь к операциям над операторами измерения. Линейная комбинация операторов измерения переводит по нашему определению объект σ в линейную комбинацию значений исходных операторов от σ . В частности, из (14) вытекает представление матрицы наблюдений

$$Y = (\hat{A} \otimes \dots \otimes \hat{\eta}) \sigma. \quad (15)$$

Введем еще векторные операторы

$$\bar{A} = \bar{U}_{n_A} \hat{A} = \begin{pmatrix} \hat{A}_1 \\ \vdots \\ \hat{A}_n \end{pmatrix}, \dots, \bar{\eta} = \bar{U}_n \hat{\eta} = \begin{pmatrix} \hat{\eta}_1 \\ \vdots \\ \hat{\eta}_n \end{pmatrix},$$

где

$$\hat{A}_1 = \frac{1}{n_A} \sum_{i=1}^{n_A} \hat{A}_i, \dots, \hat{\eta}_1 = \frac{1}{n} \sum_{y=1}^n \hat{\eta}_y.$$

Тогда имеем, например,

$$(\hat{A}_1 \dots \hat{E}_m \hat{\eta}_1) \sigma = n^{-1} \sum_{y=1}^n (\hat{A}_1 \dots \hat{E}_m \hat{\eta}_y) \sigma = (y_{1\dots m}^{(1)}, \dots, y_{1\dots m}^{(q)}), \quad (16)$$

где в правой части использована известная [1] в ДА операция пунктирования индекса (усреднение по индексу). Матрица из строк (16) имеет вид $(\hat{A} \otimes \dots \otimes \hat{E} \otimes \bar{\eta}) \sigma$.

3.10. Переходим к вычислению (10). Используя (15), немедленно получим равенство

$$y^T(Q_t - Q)y = \sigma^T(\hat{A}^T_{t_1} \hat{A} \otimes \dots \otimes \hat{A}^T_{t_{k+1}} \hat{A}) \sigma, \quad (17)$$

где (согласно 3.7) t_1 равно 1, $1-\bar{0}$ или $\bar{0}$, в зависимости от того, является ли 1-ый индекс дополнительным, коренным или отсутствующим в ψ_t . Например,

$$\hat{A}^T_{t_1} \hat{A} = \begin{cases} \sum_1 \hat{A}^T_1 \hat{A}_1, & \text{если } 1 - \text{дополнительный,} \\ \sum_1 \hat{A}^T_1 \hat{A}_., & \text{если } 1 \text{ не используется в } \psi_t, \\ \sum_1 (\hat{A}^T_1 \hat{A}_1 - \hat{A}^T_1 \hat{A}_.), & \text{если } 1 - \text{коренной.} \end{cases} \quad (18)$$

Применяя правило (18) для каждого t_1, \dots, t_{k+1} в (17), находим окончательный вид произведения (10). Например, для эффекта $\psi_t = (bd_{j1(1k)})$ в (5) получим

$$\begin{aligned} y^T(Q_t - Q)y = & \sum_{1, \dots, m, v} \sigma^T[\hat{A}^T_1 \hat{A}_1 \otimes (\hat{B}^T_j \hat{B}_j - \hat{B}^T_j \hat{B}_.) \otimes \hat{C}^T_k \hat{C}_k \otimes (\hat{D}^T_1 \hat{D}_1 - \hat{D}^T_1 \hat{D}_.) \otimes \\ & \otimes \hat{E}^T_1 \hat{E}_. \otimes \hat{\eta}^T_1 \hat{\eta}_.] \sigma = \sum_{1, \dots, m, v} (y^T_{ijk1..} y_{ijk1..} - y^T_{ijk...} y_{ijk...} - \\ & - y^T_{i.k1..} y_{i.k1..} + y^T_{i.k...} y_{i.k...}), \end{aligned} \quad (19)$$

где y с индексами обозначает строку q чисел. Напомним, что (19) выражает прирост остаточной суммы квадратов от игнорирования эффекта ψ_t в модели (2). Формула (19) непосредственно пригодна для вычисления этого прироста.

3.11. Резюмируя сделанное отметим, что ДА любого сбалансированного ПФК можно выполнить стандартной методикой, которая хорошо известна для разных частных случаев ПФК, например в [7]. При доказательстве правил ДА ПФК удобно поль-

зоваться канонической формой линейной модели и операторами измерения.

Сформулируем, наконец, основные технические приемы ДА ПФК в виде теоремы.

Теорема (дисперсионного анализа сбалансированных ПФК). В случае произвольного k -факторного сбалансированного ПФК полная квадратная сумма $Y^T Y$ ($= (q \times q)$ -матрица) разлагается в квадратные суммы отдельных эффектов и остаточную квадратную сумму $Y^T Q Y$. Эффекты соответствуют общему среднему μ и всевозможным наборам $k \geq 1$ взаимно не подчиняющихся факторов. Квадратную сумму любого эффекта можно вычислить суммированием по всем $k+1$ индексам выражения, которое является суммой всевозможных матриц вида $(-1)^{\rho} Y_{\underline{\alpha}}^T Y_{\underline{\alpha}}$. Здесь знак $\underline{\alpha}$ обозначает $(k+1)$ -местное индексное поле, где точками заменены все индексы, которые в эффекте отсутствуют, и еще $\rho \geq 0$ коренных индексов, а все дополнительные индексы эффекта обязательно указаны.

Число степеней свободы квадратной суммы эффекта равно произведению $k+1$ сомножителей, где каждому пропущенному в эффекте индексу соответствует сомножитель 1, каждому дополнительному индексу соответствует сомножитель p , а каждому коренному индексу эффекта соответствует сомножитель $p-1$ (p обозначает сокращенную размерность индекса).

Доказательство в основном уже дано. Разложение суммы $Y^T Y$ вытекает из равенства

$$Y^T Y = Y^T Q Y + \sum_{t=1}^T Y^T (Q_t - Q) Y ,$$

которое следует из (9). Число степеней свободы квадратной суммы легко вычислить как $\text{sp}(Q_t - Q)$, учтя что след тензорного произведения матриц (например, произведения (13)) равен произведению следов сомножителей.

Л и т е р а т у р а

1. Ahrens, H., *Varianzanalyse*. Berlin, 1967.
2. Möls, T., *Matemaatika bioloogidele I*. Tartu, 1978.
3. Длин А.М., *Факторный анализ в производстве*. М., 1975.
4. Кендалл М., Стьюарт А., *Многомерный статистический анализ и временные ряды*. М., 1976.
5. Рао С.Р., *Линейные статистические методы и их применения*. М., 1968.
6. Хикс Ч., *Основные принципы планирования эксперимента*. М., 1967.
7. Шеффе Г., *Дисперсионный анализ*. М., 1963.

НОВОЕ СЕМЕЙСТВО ПОКАЗАТЕЛЕЙ СТАТИСТИЧЕСКОЙ ЗАВИСИМОСТИ И ЕГО ПРИМЕНЕНИЯ

Ю. Вардья, Т. Мелс

В работе введено семейство показателей тесноты взаимной статистической зависимости двух случайных величин на конечном множестве и указаны некоторые применения этих показателей. Доказано, что рассматриваемые показатели являются корреляционными функционалами в смысле [1]. То, что они имеют тривиальную характерную группу, позволяет применять их для проверки независимости против зависимости заданного вида. Методом Монте-Карло моделирования показано, что для определенных альтернатив соответствующий критерий имеет большую мощность чем критерий χ^2 .

1. Определения и соглашения

В данной работе рассматриваются случайные величины на конечном множестве $M = \{a_1, \dots, a_n\}$. Класс всех таких случайных величин обозначим X , а подкласс невырожденных случайных величин в нем обозначим Y .

Следуя [1] будем говорить, что отображение $k: Y \times Y \rightarrow [0, 1]$ является корреляционным функционалом (КФ), если оно удовлетворяет следующим условиям 1^0-8^0 (аксиомам КФ): (1^0) существ-

вует такая группа G преобразований $M \rightarrow M$, что $k(x, y) = 1$ тогда и только тогда, когда $P\{x = gy\} = 1$ при некотором $g \in G$ (G называется характерной группой); (2°) значение $k(x, y)$ зависит от x и y через их совместное распределение $P_{x, y}$, т.е. $k(x, y) = k(P_{x, y})$; (3°) $k(x, y) = k(y, x)$; (4°) если $k(x_1, x) \rightarrow 1$, то $k(x_1, y) \rightarrow k(x, y)$ при любом $y \in Y$; (5°) если x и y статистически независимы, то $k(x, y) = 0$; (6°) если $k(x_1, x) \rightarrow 1$, то можно найти g_1, g_2, \dots в G таким образом, чтобы $g_1 x_1 \rightarrow x$ по вероятности; (7°) если x и z независимы и $0 \leq t \leq 1$, то

$$k(P_{x, y}) \geq k(tP_{x, y} + (1 - t)P_{x, z}); \quad (1)$$

(8°) при $t \rightarrow 1$ неравенство (1) переходит в равенство.

КФ называется чувствительным, если для него верно обращение аксиомы 5°, и широким, если из сходимости по вероятности $x_1 \rightarrow x$ следует сходимость $k(x_1, x) \rightarrow 1$. Известно [2], что для случайных величин на конечном множестве существуют только широкие КФ, а для любой группы преобразований $M \rightarrow M$ существует чувствительный КФ, имеющий эту группу своей характерной группой. Однако, до сих пор мало описано применений введенных понятий. В следующем пункте статьи мы вводим одно семейство чувствительных широких КФ, имеющих тривиальную характерную группу, а затем (в п. 3) описываем некоторые возможные применения этого семейства.

2. Семейство метрических корреляционных функционалов

Пусть $M = \{a_1, \dots, a_n\}$. Дадим функционал $k_{\alpha, \sigma, \nu, \pi}(x, y)$ для всех случайных величин x и y из множества Y невырожденных случайных величин правилом

$$k_{\alpha, \sigma, \nu, \kappa}(x, y) = \frac{d_{\alpha}^{\nu}(P_{x, y}; P_{x, y}^{\otimes})}{d_{\alpha}^{\nu}(P_{x, y}; P_{x, y}^{\otimes}) + \sigma d_{\alpha}^{\kappa}(P_{x, y}; P_{x, y}^{\Delta})}, \quad (2)$$

где $P_{x, y}$ — распределение вектора (x, y) , $P_{x, y}^{\otimes} = P_x^{\otimes} P_y$, d_{α} — метрика вида

$$d_{\alpha}(P, Q) = \sqrt{\sum_{i, j=1}^n \alpha_{ij} (p_{ij} - q_{ij})^2}, \quad (3)$$

($\alpha_{ij} = \alpha_{ji} > 0$, $i, j = 1, \dots, n$), $P_{x, y}^{\Delta}$ — ближайшее в метрике d_{α} к $P_{x, y}$ распределение класса $\Delta_2 = \{P: p_{ij} = 0, \text{ если } i \neq j\}$, σ — константа. В частном случае, где $\alpha_{ij} = 1$ для всех $i, j = 1, \dots, n$, сокращаем обозначение $k_{\alpha, \sigma, \nu, \kappa}$ до $k_{\sigma, \nu, \kappa}$.

Отметим, что определение функционала (2) корректно, поскольку

$$d_{\alpha}^{\nu}(P_{x, y}; P_{x, y}^{\otimes}) + \sigma d_{\alpha}^{\kappa}(P_{x, y}; P_{x, y}^{\Delta}) = 0$$

тогда и только тогда, когда $P_{x, y} = P_{x, y}^{\otimes} = P_{x, y}^{\Delta}$, откуда следует, что пара (x, y) не входит в $Y \times Y$, т.е. в область определения функционала (2).

Теорема. Функционал $k_{\alpha, \sigma, \nu, \kappa}(x, y)$ является при $0 < \kappa \leq \nu$ и $\sigma > 0$ широким чувствительным корреляционным функционалом с тривиальной характерной группой.

Доказательство. Покажем сначала, что функционал $k_{\alpha, \sigma, \nu, \kappa}$ удовлетворяет аксиомам корреляционного функционала. Удовлетворенность аксиом 1⁰, 2⁰ и 5⁰ видно непосредственно. Выполненность аксиомы 3⁰ легко вытекает из (3) и из симметрии $\alpha_{ij} = \alpha_{ji}$. Проверим аксиомы 4⁰ и 6⁰.

Если $k_{\alpha, \sigma, \nu, \kappa}(x_1, x) \rightarrow 1$, то $\sigma d_{\alpha}^{\kappa}(P_{x_1, x}; P_{x_1, x}^{\Delta}) \rightarrow 0$ или $d_{\alpha}(P_{x_1, x}; P_{x_1, x}^{\Delta}) \rightarrow 0$. Это возможно только при $P\{x_1 = x\} \rightarrow 1$

или $P_{x_1, x} \rightarrow P_{x, x}$, откуда, очевидно, следует сходимость по вероятности $x_1 \rightarrow x$. Тем самым $k_{\alpha, \sigma, \nu, \kappa}(x, y)$ удовлетворяет аксиоме 6⁰. Далее, из сходимости $x_1 \rightarrow x$ следует сходимость $P_{x_1, y} \rightarrow P_{x, y}$ и в силу непрерывности диагонализации $P \rightarrow P^\Delta$ и непрерывности метрики d_α имеем

$$d_\alpha(P_{x_1, y}; P_{x_1, y}^\Delta) \rightarrow d_\alpha(P_{x, y}; P_{x, y}^\Delta).$$

Аналогично, из непрерывности рандомизации $P \rightarrow P^\oplus$ следует сходимость

$$d_\alpha(P_{x_1, y}; P_{x_1, y}^\oplus) \rightarrow d_\alpha(P_{x, y}; P_{x, y}^\oplus).$$

Таким образом,

$$k_{\alpha, \sigma, \nu, \kappa}(x_1, y) \rightarrow k_{\alpha, \sigma, \nu, \kappa}(x, y),$$

т.е. $k_{\alpha, \sigma, \nu, \kappa}$ удовлетворяет аксиоме 4⁰.

Проверим аксиомы 7⁰ и 8⁰. Покажем сначала, что если x и z независимы, $0 < \kappa \leq \nu$ и $0 \leq t \leq 1$, то имеют место равенство

$$d_\alpha^\nu((tP_{x, y} + (1-t)P_{x, z}); (tP_{x, y} + (1-t)P_{x, z})^\oplus) = t^\nu d_\alpha^\nu(P_{x, y}; P_{x, y}^\oplus) \quad (4)$$

и неравенство

$$d_\alpha^\kappa((tP_{x, y} + (1-t)P_{x, z}); (tP_{x, y} + (1-t)P_{x, z})^\Delta) \geq t^\kappa d_\alpha^\kappa(P_{x, y}; P_{x, y}^\Delta). \quad (5)$$

Пусть распределения случайных величин (x, y) , x , y и z задаются соответственно матрицами (p_{ij}) , (p_i) , (q_j) и (s_j) . Тогда

$$\begin{aligned} d_\alpha^2((tP_{x, y} + (1-t)P_{x, z}); (tP_{x, y} + (1-t)P_{x, z})^\oplus) = \\ = \sum_{i, j=1}^n \alpha_{ij} (tp_{ij} - tp_i q_j)^2 = t^2 \sum_{i, j=1}^n \alpha_{ij} (p_{ij} - p_i q_j)^2 = t^2 d_\alpha^2(P_{x, y}; P_{x, y}^\oplus), \end{aligned}$$

откуда и вытекает (4). Пользуясь далее формулой

$$\delta_k = p_{kk} + \sum_{i \neq j} p_{ij} / \alpha_{kk} \sum_i \alpha_{ii}^{-1},$$

где δ_k - диагональный элемент матрицы $P_{x,y}^\Delta$, получим

$$\begin{aligned} d_\alpha^2(p_{x,y}; p_{x,y}^\Delta) &= \sum_k \alpha_{kk} \frac{(\sum_{i \neq j} p_{ij})^2}{\alpha_{kk}^2 (\sum_i \alpha_{ii}^{-1})^2} + \sum_{i \neq j} \alpha_{ij} p_{ij}^2 = \\ &= (\sum_{i \neq j} p_{ij})^2 / \Lambda + \sum_{i \neq j} \alpha_{ij} p_{ij}^2, \end{aligned}$$

где $\Lambda = \sum_i \alpha_{ii}^{-1}$. Следовательно,

$$\begin{aligned} d_\alpha^2((tP_{x,y} + (1-t)P_{x,z}); (tP_{x,y} + (1-t)P_{x,z})^\Delta) &= \\ &= (\sum_{i \neq j} (tp_{ij} + (1-t)p_{iz}))^2 / \Lambda + \sum_{i \neq j} \alpha_{ij} (tp_{ij} + (1-t)p_{iz})^2 \geq \\ &\geq t^2 (\sum_{i \neq j} p_{ij})^2 / \Lambda + t^2 \sum_{i \neq j} \alpha_{ij} p_{ij}^2 = t^2 d_\alpha^2(p_{x,y}; p_{x,y}^\Delta), \end{aligned}$$

что и доказывает (5).

Возвращаемся к доказательству аксиом 7^0 и 8^0 . Пользуясь формулами (2) и (4) и неравенством $0 \leq t \leq 1$, получим

$$k_{\alpha, \rho, \nu, \kappa}(P_{x,y}) \frac{d_\alpha^\nu((tP_{x,y} + (1-t)P_{x,z}); (tP_{x,y} + (1-t)P_{x,z})^\otimes)}{d_\alpha^\nu((tP_{x,y} + (1-t)P_{x,z}); (tP_{x,y} + (1-t)P_{x,z})^\otimes) + t^\nu d_\alpha^\kappa(P_{x,y}; P_{x,y}^\Delta)}.$$

Учтя неравенство $t^\nu \leq t^\kappa$ и формулу (5) увидим, что действительно

$$k_{\alpha, \rho, \nu, \kappa}(P_{x,y}) \geq k_{\alpha, \rho, \nu, \kappa}(tP_{x,y} + (1-t)P_{x,z}). \quad (6)$$

Таким образом, аксиома 7^0 выполнена. Легко видеть, что при $t \rightarrow 1$ неравенство (6) превратится в равенство и поэтому

$k_{\alpha, \sigma, \nu, \kappa}$ удовлетворяет также и аксиоме 8⁰.

Мы уже доказали, что $k_{\alpha, \sigma, \nu, \kappa}$ является КФ. Осталось установить его чувствительность, так как широкость вытекает из общей теории КФ. Допустим, что x и y статистически зависимы. Тогда $d_{\alpha}^{\nu}(P_{x,y}; P_{x,y}^{\otimes}) > 0$ и, следовательно, $k_{\alpha, \sigma, \nu, \kappa}(x, y) > 0$. Значит, $k_{\alpha, \sigma, \nu, \kappa}$ — чувствительный КФ. Теорема полностью доказана.

КФ, определенный формулой (2), будем в дальнейшем называть метрическим КФ. Укажем здесь некоторые простейшие свойства метрических КФ (в их истинности легко убедиться, представив КФ (2) в виде

$$k_{\alpha, \sigma, \nu, \kappa}(x, y) = \frac{1}{1 + \text{od}_{\alpha}^{\kappa}(P_{x,y}; P_{x,y}^{\Delta}) / d_{\alpha}^{\nu}(P_{x,y}; P_{x,y}^{\otimes})} :$$

при $\sigma_1 \geq \sigma_2$

$$k_{\alpha, \sigma_1, \nu, \kappa}(x, y) \leq k_{\alpha, \sigma_2, \nu, \kappa}(x, y);$$

при $\nu_1 \geq \nu_2$

$$k_{\alpha, \sigma, \nu_1, \kappa}(x, y) \geq k_{\alpha, \sigma, \nu_2, \kappa}(x, y);$$

при $\kappa_1 \geq \kappa_2$

$$k_{\alpha, \sigma, \nu, \kappa_1}(x, y) \leq k_{\alpha, \sigma, \nu, \kappa_2}(x, y);$$

при $\kappa_1/\nu_1 = \kappa_2/\nu_2$

$$k_{\alpha, \sigma_2, \nu_2, \kappa_2}(x, y) = (1 + \sigma_2((k_{\alpha, \sigma_1, \nu_1, \kappa_1}^{-1}(x, y) - 1)/\sigma_1)^{\nu_2/\nu_1 - 1});$$

при $\alpha = b \cdot \alpha'$

$$k_{\alpha, \sigma, \nu, \kappa}(x, y) = k_{\alpha', \sigma \cdot b^{\frac{\nu - \nu'}{\nu}}, \nu, \kappa}(x, y).$$

3. Применение метрических корреляционных функционалов к исследованию статистической зависимости

Рассматриваем канал связи с n символами a_1, \dots, a_n , входом x и выходом y . Статистическую зависимость y от x интерпретируем как передачу информации через канал. Однако связь устроена так, что выходной символ принимается как входной, т.е. преобразования символов в канале не предполагается. Поэтому передаваемая (с помощью статистической зависимости) информация может иметь практическую ценность лишь при условии, что вероятность совпадения входных и выходных символов в некотором среднем смысле превышает вероятность их случайного совпадения. В таком случае будем говорить, что имеется прямая зависимость между входом и выходом.

При изучении описанного канала иногда возникает необходимость измерять силу прямой статистической зависимости между входом и выходом, а также проверить реальность этой зависимости. Мы предлагаем использовать для решения этих или аналогичных проблем КФ с тривиальной характерной группой, поскольку они более чувствительны к прямой зависимости аргументов.

Некоторые метрические КФ также имеют достаточно малую чувствительность к непрямой зависимости аргументов. Чтобы иллюстрировать возможности применения метрических КФ, рассматриваем ради простоты двоичный канал с передаваемыми символами a_1 и a_2 . Если маргинальные распределения P_x и P_y фиксированы для такого канала, то вероятность $p_{11} = P\{x = a_1 \text{ и } y = a_2\}$ однозначно определяет совместное распределение $P_{x,y}$.

Поэтому в данном случае в качестве аргумента для КФ достаточно взять p_{11} .

На рисунке 1 показаны графики некоторых метрических КФ как функции от p_{11} . Эти графики показывают, как можно подобрать при фиксированных P_x и P_y параметры α и β , чтобы значение $k_{\alpha, \beta}$ превышал бы заданный уровень (например, 0.09) лишь при достаточно сильной прямой зависимости (т.е. при $p_{11} \gg p_{1 \cdot} p_{\cdot 1}$, где $p_{1 \cdot} = P\{x = a_1\}$ и $p_{\cdot 1} = P\{y = a_1\}$), и в этом случае значение $k_{\alpha, \beta}$ близко к 1.

Теперь допустим, что распределение $P_{x,y}$ оценивается эмпирическим распределением (которое обозначаем $P_{x,y}^B$), а $k_{\alpha, \beta}(P_{x,y}^B)$ служит статистиком для проверки прямой зависимости между x и y . Тогда область нечувствительности $k_{\alpha, \beta}$ охватывает и некоторую долю случаев (несильной) прямой статистической зависимости (см. рис. 1). Подбором параметров α и β эта область может быть совмещена с областью принятия нулевой гипотезы H_0 об отсутствии прямой зависимости. В таком случае нулевая гипотеза принимается при $k_{\alpha, \beta}(P_{x,y}^B) \approx 0$ и отклоняется при $k_{\alpha, \beta}(P_{x,y}^B) \approx 1$. Иными словами, статистик $k_{\alpha, \beta}(P_{x,y}^B)$ является своего рода индикатором критической области этой гипотезы.

Применение описанного критерия было статистически моделировано (на микрокалькуляторе TI-58). По результатам этого моделирования построено эмпирическое распределение величины $k_{0.01, 10, 10}(P_{x,y}^B)$, где каждое эмпирическое распределение $P_{x,y}^B$ получено по $n = 15$ измерениям — имитациям вектора (x, y) с маргинальными распределениями $P_x = P_y = (0.5, 0.5)$ и заданным значением p_{11} . Сводка результатов этого моделирования

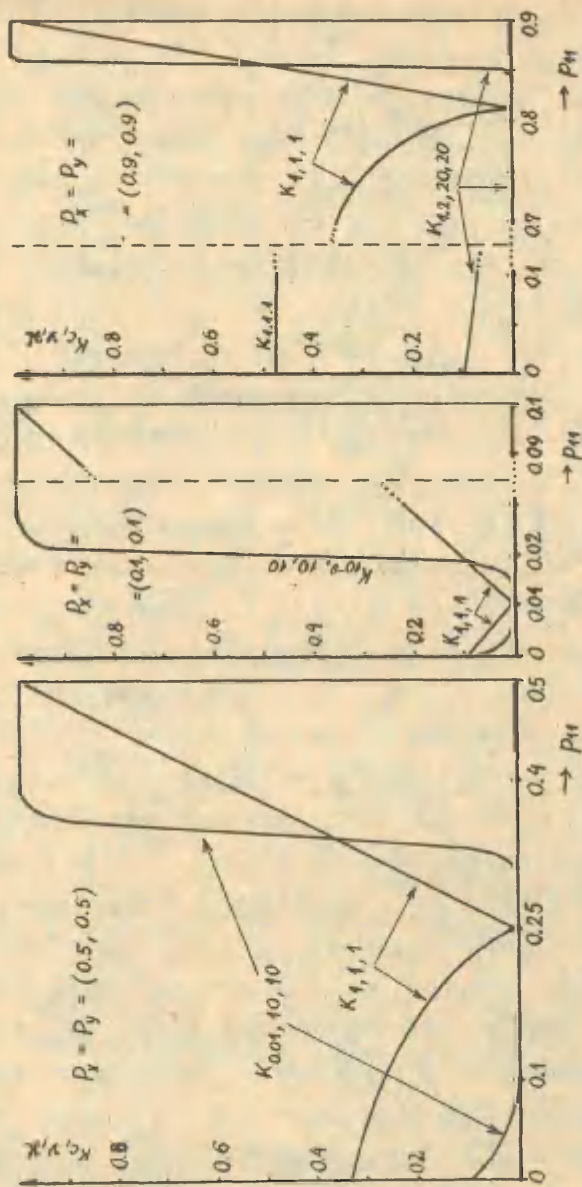


Рис. 1. Графики некоторых метрических корреляционных функционалов в двумерном случае.

Таблица 1.

Сводные результаты статистического моделирования применений

критериев $k_{0.01, 10, 10}(P_{x,y}^a)$ и χ^2 при $P_{x,y} = (0.5, 0.5)$ и $n=15$.В скобках даны результаты, касающиеся критерия χ^2 .

Характер статистической зависимости	P_{11}	Эмпирич. распредел. $k_{0.01, 10, 10}(P_{x,y}^a)$						Число моделир. $P_{x,y}^a$	Эмпирич. уровень значимости	Эмпирич. мощность критерия
		[0.0, 0.1)	[0.1, 0.2)	[0.2, 0.5)	[0.5, 0.8)	[0.8, 0.9)	[0.9, 1.0]			
непрямая зависимость	0.10							(105)	(0.63)	
	0.15							(121)	(0.32)	
	0.20	0.98	0.00	0.00	0.01	0.01	0.00	152	0.02	
независимость	0.25	0.94	0.00	0.00	0.01	0.03	0.00	(116)	(0.09)	
								259	0.058	
	0.30	0.75	0.01	0.00	0.07	0.17	0.00	(211)	(0.062)	0.25
прямая зависимость	0.35	0.48	0.01	0.00	0.12	0.39	0.01	(110)		(0.15)
								124		0.52
	0.40	0.16	0.00	0.00	0.04	0.79	0.02	(150)		(0.30)
								107		0.84
	0.45	0.02	0.00	0.00	0.03	0.80	0.14	(116)		(0.68)
	0.50							209		0.98
								(142)		(0.97)
								-		1.00

приведена в таблице 1.

Как видно, в данном случае значения $k_{0.01,10,10}$ редко попадают в интервал $[0.1, 0.8)$, что и следовало бы ожидать (см. рис. 1). Из таблицы 1 и построенного по ней рисунка 2 вытекает также, что для рассматриваемой пары конкурирующих гипотез мощность критерия, построенного по $k_{0.01,10,10}(P_{x,y}^3)$, превышает мощность обычного критерия χ^2 .

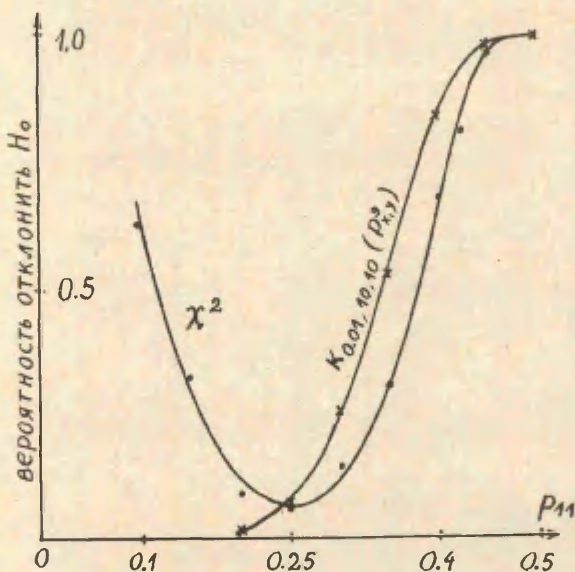


Рис. 2. Кривые мощности двух критериев.

Л и т е р а т у р а

1. Мелс Т.Э., Определение общих корреляционных коэффициентов. Труды ВЦ ТТУ, 1972, 25, 64-78.
2. Мелс Т.Э., Описание класса общих корреляционных коэффициентов для случайных элементов на конечном множестве. Труды ВЦ ТТУ, 1972, 26, 3-18.

АСИМПТОТИЧЕСКОЕ РАСПРЕДЕЛЕНИЕ
КОЭФФИЦИЕНТОВ ПОЛИНОМИАЛЬНОЙ РЕГРЕССИИ

А.-М. Парринг

1. Введение

Пусть дан случайный вектор $x'_* = (x_0, x_1, \dots, x_p)$. Функцией регрессии для компонента x_0 называется функция $g(x)$, выполняющая условие

$$E(x_0 - g(x))^2 = \min_{\varphi \in G} (x_0 - \varphi(x))^2,$$

где G — класс измеримых функций. Вычисление (оценивание) функции регрессии возможно в случае, когда известно распределение вектора x_* (распределение известно до точности неизвестных параметров, т.е. возможно определить аналитический вид функции $g(x)$). Меньше априорной информации требует определение наилучшего приближения к компоненту x_0 из некоторого класса функции \mathcal{F} , т.е. функции $h(x)$, которая удовлетворяет условию

$$E(x_0 - h(x))^2 = \min_{\varphi \in \mathcal{F}} E(x_0 - \varphi(x))^2.$$

Функцию $h(x)$ будем называть функцией регрессии из класса \mathcal{F} . Как известно (см. [3], стр. 126), функция $h(x)$ является наилучшим приближением функции регрессии из класса \mathcal{F} в смысле

наименьших квадратов.

В дальнейшем посмотрим подробнее случай, где классом \mathcal{F} выбран класс многочленов степени k из компонентов вектора x .

2. Функция полиномиальной регрессии

Пусть классом $\mathcal{F}_{k,I}$ является класс многочленов степени k из компонентов вектора x , т.е. если $f \in \mathcal{F}_{k,I}$, то

$$f(x) = \alpha_0 + \sum_{(l_1, \dots, l_p) \in I} \alpha_{l_1 \dots l_p} x_1^{l_1} \dots x_p^{l_p},$$

где $I = \{(l_1, \dots, l_p) : 0 \leq l_s \leq k, \sum_{s=1}^p l_s \leq k\}$ — некоторое фиксированное множество. Например, если $p=2$ и

$$I = \{(1,0), (0,1), (1,1)\}, \quad (1)$$

то $f(x) = \alpha_0 + \alpha_{10}x_1 + \alpha_{01}x_2 + \alpha_{11}x_1x_2$. Функцию регрессии из класса $\mathcal{F}_{k,I}$ называем функцией полиномиальной регрессии.

Пусть множество I имеет s элементов. Обозначим вектор, компонентами которого являются упорядоченные некоторым образом произведения $x_1^{l_1} \dots x_p^{l_p}$ через z , $z \in \mathcal{M}^{s \times 1}$, а вектор, компонентами которого являются таким же образом упорядоченные коэффициенты $\alpha_{l_1 \dots l_p}$, через x . Тогда функция полиномиальной регрессии имеет вид $f(x) = x_0 + x'z$. Предположим, что среди компонентов вектора x нет линейно зависимых и ни один из компонентов вектора x не является полиномиальной функцией от остальных. Так как функция полиномиальной регрессии явля-

ется линейной функцией вектора z , то вектор коэффициентов² $x_* = [x_0 : x']$ определяется системой линейных уравнений (см. напр. [3], стр. 129)

$$\begin{cases} x = \Sigma^{-1} \epsilon_0 \\ x_0 = \mu_0 - x' \mu, \end{cases} \quad (2)$$

где через $\Sigma_* = \begin{bmatrix} \epsilon_{00} & \epsilon'_0 \\ \epsilon_0 & \Sigma - 1 \end{bmatrix}$ обозначена дисперсионная матрица, а через $\mu_* = [\mu_0 : \mu']$ вектор средних вектора z_* , $z'_* = [z_0 : z']$, $z_0 = x_0$. Тесноту связи измеряет множественный коэффициент корреляции (МКК)

$$\rho^2 = \epsilon'_0 \Sigma^{-1} \epsilon_0 / \epsilon_{00}. \quad (3)$$

Выборочную оценку a_* вектора x_* получим, заменив в системе (2) дисперсионную матрицу и вектор средних их выборочными оценками S_* и \bar{x}_* . Асимптотическое распределение вектора оценок изучено в [2]. Оказывается, что независимо от распределения вектора z_* , асимптотическим распределением вектора оценок является нормальное распределение. Параметры этого распределения зависят от центральных моментов третьего и четвертого порядка.

3. Асимптотическое распределение коэффициентов полиномиальной регрессии

Введем следующую символику. Обозначим через $*L_k$ случайную матрицу, определенную рекурсивным соотношением³

² Квадратные скобки указывают блочную структуру матрицы. Так блок $x' = (x_1, \dots, x_s)$.

$$*L_k = \begin{cases} *L_{k-1} \otimes z_*', & \text{если } k \text{ четное} \\ *L_{k-1} \otimes z_*, & \text{если } k \text{ нечетное} \end{cases}$$

и через $*\bar{L}_k$ аналогичную матрицу, определенную для центрированного случайного вектора $z_* - \mu_*$

$$*\bar{L}_k = \begin{cases} *\bar{L}_{k-1} \otimes (z_* - \mu_*)', & \text{если } k \text{ четное} \\ *\bar{L}_{k-1} \otimes (z_* - \mu_*), & \text{если } k \text{ нечетное,} \end{cases}$$

$$k=1, 2, \dots; \quad *L_0 = *\bar{L}_0 = 1.$$

Матрицу $*M_k = E_* L_k$ называем матрицей k -моментов случайного вектора z_* , матрицу $*\bar{M}_k = E_* \bar{L}_k$ — матрицей центральных k -моментов случайного вектора z_* . Если все математические ожидания, являющиеся элементами матрицы $*M_k$, существуют и конечны, то скажем, что вектор z_* имеет матрицу k -моментов. Кратко напомним в этом случае $z_* \in \mathcal{H}^k$.

Таким образом, матрицы $*\bar{M}_3$ и $*\bar{M}_4$ являются матрицами третьих и четвертых центральных моментов вектора z_* . В дальнейшем используем следующее расчленение этих матриц в блоки: $*\bar{M}_4 = [(*\bar{M}_4)_{ij}]$, $i, j=0, 1, \dots, s$, где блок $(*\bar{M}_4)_{ij} = E((x_i - \mu_1)(x_j - \mu_1) * \bar{L}_2)$ и $*\bar{M}_3 = [(*\bar{M}_3)_i]$, $i=0, 1, \dots, s$, где блок $(*\bar{M}_3)_i = E((x_i - \mu_1) * \bar{L}_2)$.

Параметры асимптотического распределения вектора оценок a_* определяются следующей теоремой, доказанной в [2].

Теорема 1. Пусть случайный вектор $z_* \in \mathcal{H}_4$ и среди компонентов вектора z_* нет линейно зависимых. Тогда⁴

³ Символом \otimes обозначено прямое произведение матриц, т.е. если $A \in \mathbb{M}^{m \times n}$, то $A \otimes B = [a_{ij} B]$, $i=1, \dots, m$; $j=1, \dots, n$.

⁴ Символом $x \xrightarrow{d} y$ обозначается сходимость по распределению.

$$\sqrt{n} (a_{**} - \alpha_{**}) \xrightarrow{d} N(0, \xi),$$

где

$$\xi = \begin{bmatrix} \sigma_{00}(1-\varphi^2) - 2\mu' \Sigma^{-1} f + \mu' \Pi \mu & f' \Sigma^{-1} - \mu' \Pi \\ \Sigma^{-1} f - \Pi \mu & \Pi \end{bmatrix}, \quad (4)$$

а $\Pi = \Sigma^{-1} C \Sigma^{-1}$. Элементы c_{ij} матрицы $C \in \mathbb{M}^{s \times s}$ определены равенством $c_{ij} = \gamma' ({}_{**}\bar{M}_4)_{ij} \gamma$, где $\gamma' = [1 : -\alpha]$. Элементы f_1 вектора $f \in \mathbb{M}^{s \times 1}$ определены равенством $f_1 = \gamma' ({}_{**}\bar{M}_3)_1 \gamma$.

В случае полиномиальной регрессии не все элементы матриц ${}_{**}\bar{M}_4$ и ${}_{**}\bar{M}_3$ являются центральными моментами первоначального вектора x_{**} . Например, если множество I определено равенством (1), т.е. $z = (x_1, x_2, x_1 x_2)$, то элемент 1,2 блока $({}_{**}\bar{M}_4)_{33}$ $m_{3,3,1,2} = E((x_1 x_2 - E x_1 x_2)^2 (x_1 - E x_1)(x_2 - E x_2))$ и таким образом не является центральным моментом вектора $x' = (x_1, x_2)$. Но так как имеют место соотношения (см. [1])

$${}_{**}\bar{M}_3 = {}_{**}M_3 - \mu_{**} \otimes {}_{**}M_2 - {}_{**}M_2 \otimes \mu_{**} - \text{vec } {}_{**}M_2 \otimes \mu_{**}' + 2\mu_{**} \otimes \mu_{**}' \otimes \mu_{**}', \quad (5)$$

$$\begin{aligned} {}_{**}\bar{M}_4 = & {}_{**}M_4 - \mu_{**} \otimes {}_{**}M_3 - {}_{**}M_3 \otimes \mu_{**} - \mu_{**}' \otimes {}_{**}M_3 - {}_{**}M_3 \otimes \mu_{**}' + \\ & + \mu_{**} \otimes \mu_{**}' \otimes (\text{vec } {}_{**}M_2)' + \text{vec } {}_{**}M_2 \otimes \mu_{**}' \otimes \mu_{**}' + \mu_{**}' \otimes \mu_{**} \otimes {}_{**}M_2 + \\ & + \mu_{**}' \otimes {}_{**}M_2 \otimes \mu_{**} + \mu_{**} \otimes {}_{**}M_2 \otimes \mu_{**}' + {}_{**}M_2 \otimes \mu_{**} \otimes \mu_{**}' - \\ & - 3\mu_{**} \otimes \mu_{**}' \otimes \mu_{**} \otimes \mu_{**}', \end{aligned} \quad (6)$$

то элементы матриц ${}_{**}\bar{M}_4$ и ${}_{**}\bar{M}_3$ определяются через моменты первоначального вектора x_{**} . Предположения теоремы 1 выполнены, если $x_{**} \in \mathcal{H}^{2k}$.

4. Пример

Просмотрим конкретный пример использования приведенной теоремы. Пусть изучаемый случайный вектор двумерный, $x_* = (x_0, x_1)$. Компонент x_1 — случайная величина с экспоненциальным распределением, т.е.

$$f_{x_1}(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0, \end{cases} \quad \lambda > 0.$$

Компонент $x_0 = \frac{\nu + \lambda}{\lambda} e^{-\nu x_1}$, $\nu > 0$. Предположим, что связь между компонентами x_0 и x_1 неизвестно экспериментатору и информацию об этом векторе он может получить только в форме выборки, регистрируя конкретные значения вектора x_* . Он выбирает функцию регрессии из класса квадратных полиномов $f(x) = x_0 + x_1 x + x_2 x^2$ (в данном случае $I = \{1, 2\}$). Таким образом, $x_* = (x_0, x_1, x_1^2)$. Используя в формулах (2) и (3) выборочные оценки моментов, можно вычислить оценки коэффициентам полиномиальной регрессии и МКК.

Так как нам известно конкретный механизм изучаемого явления, можно вычислить истинные x_* , σ^2 и дисперсионную матрицу асимптотического распределения. Чтобы характеризовать поведение оценки a_* в данной ситуации, проделаем эти вычисления.

Из определения вектора x_* следует, что $x_* \in \mathcal{H}^k$ при всех конечных k . Имеет место соотношение

$$E x_1^s x_0^s = \frac{(\nu + \lambda)^s}{\lambda^{s-1}} \int_0^\infty x^s e^{-(\nu + \lambda)x} dx = \frac{(\nu + \lambda)^s}{\lambda^{s-1}} \frac{1!}{(\nu + \lambda)^{1+1}}. \quad (7)$$

Отсюда получим

$$\mu_{**}' = (Ex_0, Ex_1, Ex_1^2) = (1, 1/\lambda, 2/\lambda^2)$$

и

$$\Sigma_{**} = E(z_{**} - \mu_{**})(z_{**} - \mu_{**})' = \begin{pmatrix} \frac{v^2}{\lambda(2v+\lambda)} & -\frac{v}{\lambda(v+\lambda)} & -\frac{2v(v+2\lambda)}{\lambda^2(v+\lambda)^2} \\ -\frac{v}{\lambda(v+\lambda)} & \frac{1}{\lambda^2} & \frac{4}{\lambda^3} \\ -\frac{2v(v+2\lambda)}{\lambda^2(v+\lambda)^2} & \frac{4}{\lambda^3} & \frac{20}{\lambda^4} \end{pmatrix}.$$

Вектором коэффициентов квадратной регрессии является вектор $\alpha_{**} = [\alpha_0 : \alpha]$, где

$$\alpha = -\lambda v / 2(v+\lambda)^2 \begin{pmatrix} 2(3v+\lambda) \\ -\lambda v \end{pmatrix}$$

и

$$\alpha_0 = (\lambda^2 + 3v\lambda + 3v^2) / (\lambda + v)^2,$$

таким образом

$$f(z) = \frac{\lambda^2 + 2v(\lambda + v)}{(\lambda + v)^2} - \frac{2v(3v + \lambda)}{(\lambda + v)^2} x_1 + \frac{\lambda^2 v^2}{2(\lambda + v)^2} x_1^2.$$

МКК имеет значение

$$\rho^2 = \lambda(2v + \lambda)(2v^2 + 2v\lambda + \lambda^2) / (v + \lambda)^4.$$

Чтобы избежать лишних сложностей в дальнейших вычислениях, ограничимся с случаем $v = \lambda = 1$. При таких значениях параметров $\rho^2 = 15/16 = 0,9375$ — значение, которое укажет, что выбранная функция является весьма хорошим приближением к компоненту x_0 . Но посмотрим теперь дисперсионную матрицу асимптотического распределения. Вычислив при помощи формул (5), (6), (7) матрицы \bar{M}_4 и \bar{M}_3 и имея в виду, что при данных значениях параметров $\gamma' = (1, 1, -0,125)$, получим

$$\hat{\xi} = \begin{pmatrix} 3,7657 & -8,8663 & 2,4969 \\ -8,8663 & 19,7837 & -5,5895 \\ 2,4969 & -5,5895 & 1,5850 \end{pmatrix}.$$

Если объем выборки n большой, дисперсионную матрицу оценки можно читать приблизительно равной $1/n\hat{\xi}$.

Для сравнения вычислим еще дисперсионную матрицу асимптотического распределения коэффициентов регрессии вектора $y_*' = (y_0, y_1, y_2)$, первые и вторые моменты которого совпадают с соответствующими моментами вектора z_* при $\lambda = \nu = 1$, но который имеет нормальное распределение. Используя в [2] полученный результат, можем написать

$$\xi = \xi_{00}(1-\rho^2) \begin{bmatrix} 1+\mu'\Sigma^{-1}\mu & \mu'\Sigma^{-1} \\ \Sigma^{-1}\mu & \Sigma^{-1} \end{bmatrix} = \begin{pmatrix} 0,0625 & 0,0625 & 0,0104 \\ 0,0625 & 0,1040 & -0,0208 \\ 0,0104 & -0,0208 & 0,0052 \end{pmatrix}.$$

Как увидим, при том самом объеме выборки приближенная оценка $1/n\hat{\xi}$ дисперсионной матрицы коэффициентов значительно меньше, чем в предыдущем случае.

Но и при векторе x_* точность оценок резко повышается, если экспериментатор догадается провести вместо оценки функции квадратной регрессии преобразование данных и оценить для вектора $w_*' = (w_0, w_1) = (\ln x_0, w_1)$ функцию линейной регрессии. Действительно, в этом случае $w_0 = \ln \frac{\nu+1}{\lambda} - \nu w_1$, $\rho^2 = 1$, и при объеме выборки $n \geq 2$ мы получим точные значения коэффициентов функции линейной регрессии. Вычислив в этом случае дисперсионную матрицу асимптотического распределения коэффициентов, увидим, что этой матрицей является нулевая матрица. Действительно, теперь вектор $\eta' = (1, \nu)$. Так как $w_0 - Ew_0 =$

$= -v(x_1 - 1/\lambda)$, можем написать, обозначив $E(x_1 - 1/\lambda)^k$ через a^k ,

$$(w_{**} \bar{M}_4)_{1j} = (-v)^{2-1-j} a^4 \begin{pmatrix} v^2 & -v \\ -v & 1 \end{pmatrix}$$

и

$$(w_{**} \bar{M}_3)_{1j} = (-v)^{1-1} a^3 \begin{pmatrix} v^2 & -v \\ -v & 1 \end{pmatrix},$$

$1, j=0, 1$.

Следовательно

$$\gamma'(w_{**} \bar{M}_4)_{1j} \gamma = \gamma'(w_{**} \bar{M}_3)_{1j} \gamma = 0,$$

$1, j=0, 1$. Имея в виду, что $\rho^2 = 1$, все элементы дисперсионной матрицы асимптотического распределения, определенной формулой (4), равны нулю.

Л и т е р а т у р а

1. Парринг А.-М., Вычисление асимптотических характеристик функции выборки. Уч. зап. ТГУ, № 429, 86-90.
2. Парринг А.-М., Оценка функции линейной регрессии и ее асимптотическое поведение. Уч. зап. ТГУ, № 429, 94-99.
3. Tiit, E., Parring, A., Mõls, T., Tõenäosusteooria ja matemaatilise statistika. Tallinn, 1977.

СРАВНЕНИЕ РЕГРЕССИИ В РАЗНЫХ ПОДСОВОКУПНОСТЯХ

А.Левисто, Э.Тийт

0. Одной из наиболее часто применяемых методик прикладной математической статистики является линейный регрессионный анализ. Это дает модель в виде формулы

$$y \approx \hat{y} = \beta x \quad (1)$$

для прогнозирования функционального признака y или функционального признак-вектора $y = (y_1, \dots, y_p)'$ по аргумент-признак-вектору $x = (1, x_1, \dots, x_m)'$. Параметрами прогноза являются коэффициенты регрессии, т.е. элементы матрицы коэффициентов регрессии β

$$\beta = \begin{pmatrix} \beta_{10} & \beta_{11} & \dots & \beta_{1m} \\ \dots & \dots & \dots & \dots \\ \beta_{p0} & \beta_{p1} & \dots & \beta_{pm} \end{pmatrix}. \quad (2)$$

В одномерном случае матрица β заменяется вектор-строкой $\beta' = (\beta_0, \beta_1, \dots, \beta_m)$.

Довольно типичной с точки зрения анализа данных является следующая задача: Генеральная совокупность, для которой необходимо построить модель, состоит из некоторых подсовокупностей, пусть их число K . Исследователю неизвестно, является ли для всех этих подсовокупностей пригодной единая модель

(1), или необходимо для каждой подсовокупности построить отдельную модель, т.е., вычислить для каждой подсовокупности отдельную матрицу (2).

Аналогичная задача возникает и при исследованиях, проведенных в разные моменты времени: необходимо выяснить, изменяется ли существенно в течение времени модель, описывающая исследуемое явление.

В настоящей статье, исходя из общей теории наименьших квадратов (см. [1]), выведены формулы для эффективного решения поставленной задачи при предположении нормальности рассматриваемых признак-векторов.

Основная часть излагаемой методики реализована в ВЦ ТТУ на ЭВМ ЕС-1022 в виде опытного комплекта программ.

1. Пусть имеется K подсовокупностей генеральной совокупности (или K генеральных совокупностей). В каждой из них измерен признак-вектор¹ $z = (x' : y')$, состоящий из m аргумент-признаков x_1, \dots, x_m и p функциональных признаков y_1, \dots, y_p . Пусть объем k -той выборки — n_k . Обозначаем результаты измерения в k -той совокупности через две матрицы x^k и y^k ($k = 1, 2, \dots, K$)

$$x^k = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11}^k & x_{12}^k & \dots & x_{1n_k}^k \\ \dots & \dots & \dots & \dots \\ x_{m1}^k & x_{m2}^k & \dots & x_{mn_k}^k \end{pmatrix}, \quad y^k = \begin{pmatrix} y_{11}^k & y_{12}^k & \dots & y_{1n_k}^k \\ \dots & \dots & \dots & \dots \\ y_{p1}^k & y_{p2}^k & \dots & y_{pn_k}^k \end{pmatrix}.$$

¹

Знак ' обозначает транспонирование.

Если $p=1$, то матрица Y^k заменяется на соответствующий вектор $(Y^k)' = (y_1^k, \dots, y_{n_k}^k)$.

Кроме того рассматривают объединенную выборку объема n , $n = \sum_{k=1}^K n_k$. Заметим, что объемы подвыборок n_k произвольны, и тем объединенная выборка не должна быть представительным для генеральной совокупности.

Объединенная выборка описывается с помощью матриц X и Y ,

$$X = (X^1 : \dots : X^K), \quad Y = (Y^1, \dots, Y^K), \quad (3)$$

для которых применяется упрощенное обозначение

$$X = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_{11} & x_{12} & \dots & x_{1n} \\ \dots & \dots & \dots & \dots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix}, \quad Y = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ \dots & \dots & \dots & \dots \\ y_{m1} & y_{m2} & \dots & y_{mn} \end{pmatrix}.$$

Все статистики, вычисленные по k -той выборке и по объединенной выборке обозначаются соответственно с индексом k сверху или без индекса.

Пусть в k -той генеральной совокупности имеет место соотношение (1) с матрицей коэффициентов регрессии β^k ($k=1, 2, \dots, K$), а выборочная оценка матрицы β^k обозначается через B^k , $B^k = (b_{ij}^k)$ ($i=1, \dots, p$; $j=0, 1, \dots, m$),

$$B^k = Y^k(X^k)'(X^k(X^k)')^{-1}$$

(при предположении регулярности матрицы $X^k(X^k)'$).

Поставленная задача сводится к проверке нулевой гипотезы

$$H_0 : \beta^1 = \beta^2 = \dots = \beta^K.$$

2. Для проверки этой гипотезы вычислим в каждой выборке

по принципу наименьших квадратов остатки прогноза V_0^k и V_1^k . В случае $p > 1$ они являются $(p \times p)$ -квадратичными положительно определенными матрицами, на главной диагонали которых находятся ошибки прогноза отдельных компонент функционального признак-вектора. В случае $p=1$ V_0^k и V_1^k — квадратичные формы, для которых применяется традиционная символика $(V_0^k)^2$ и $(V_1^k)^2$. Имеем

$$V_0^k = (Y^k - V^k X^k)(Y^k - V^k X^k)' = Y^k(Y^k)' - Y^k(X^k)'(X^k(X^k)')^{-1}X^k(Y^k)'.$$

Если верна нулевая гипотеза, то $\beta^k = \beta$ ($k=1, 2, \dots, K$) и при регулярности матрицы XX' для β существует оценка V , $V = (v_{ij})$ ($i=1, \dots, p$; $j=0, 1, \dots, m$),

$$V = YX'(XX')^{-1}.$$

Остаток прогноза V_1^k при справедливости нулевой гипотезы вычисляется по матрице V :

$$V_1^k = (Y^k - V^k X^k)(Y^k - V^k X^k)' \quad (k=1, 2, \dots, K).$$

Отсюда найдутся и суммарные остатки прогноза, при упрощении которых пользуются определением (3) матриц X и Y :

$$V_0 = \sum_{k=1}^K V_0^k = YY' - \sum_{k=1}^K [Y^k(X^k)'(X^k(X^k)')^{-1}X^k(Y^k)'], \quad (4)$$

$$V_1 = \sum_{k=1}^K V_1^k = YY' - YX'(XX')^{-1}XY'. \quad (5)$$

Для получения критерия проверки гипотез предполагаем, что значения аргумент-признак-вектора X фиксированы, а функциональный признак-вектор Y имеет в каждой генеральной совокупности нормальное распределение, $Y^k \sim N_p(\beta^k X^k, \Sigma_Y)$.

При справедливости нулевой гипотезы признак-вектор \mathcal{Y} имеет и в объединенной генеральной совокупности такое же нормальное распределение, $\mathcal{Y} \sim N_p(\beta X, \Sigma)$.

3. Если функциональный признак-вектор \mathcal{Y} одномерен, то задача проверки гипотез точно решается при помощи F -распределения. Действительно, в таком случае отношение

$$\frac{R_1^2 - R_0^2}{R_0^2} \cdot \frac{n - r}{w},$$

где вместо R_1 и R_0 для сохранения традиционной записи написаны R_1^2 и R_0^2 , имеет F -распределение с числами степеней свободы $n-r$ и w , где r - ранг матрицы X , n - число наблюдений и w - число независимых параметров, входящих в проверяемую гипотезу. В рассматриваемом случае $r = m+1$, $w = (K-1)(m+1)$.

4. Если же $p > 1$, то следует рассмотреть отношение детерминантов матриц остатков прогноза R_0 и R_1 . Известно (см. [1]), что отношение

$$|R_0| / |R_1|$$

имеет т.н. Λ -распределение, характеризуемое параметрами p , m и n . Для практических целей применимо приближение Λ -распределения через F :

$$\frac{1 - \Lambda^{1/s}}{\Lambda^{1/s}} \cdot \frac{ws - 2v}{mp} \sim F(pm, ws - 2v),$$

где $w = n - (p+m+1)/2$, $s = [(p^2 m^2 - 4)/(p^2 + m^2 - 5)]^{1/2}$, $v = (mp - 2)/4$; при нецелочисленных значениях величины $ws - 2v$ пользуются округленными до целых чисел значениями.

5. Иногда представляет интерес проверить гипотезу, согласно которой в разных подсовокупностях (в разные моменты времени) модели совпадают до точности свободного члена.

Для формулировки этой гипотезы вводим следующие обозначения: $\beta = (\beta_0 : \beta^*)$, $b = (b_0 : b^*)$, $b^k = (b_0^k : b^{k*})$, где $\beta_0 = (\beta_{10}, \dots, \beta_{p0})'$, $\beta^* = (\beta_{1j})$ ($i=1, \dots, p$; $j=1, \dots, m$), а b, b^k , b^* и b^{k*} — оценки матриц β и β^* соответственно по объединенной и k -той выборке. Аналогично, пусть $x' = (1, x^{*'})$, $x^* = (x_1, \dots, x_m)'$ и соответственно выборки разделяются в блоки: $x = \begin{pmatrix} 1 \\ x^* \end{pmatrix}$, $x^k = \begin{pmatrix} 1 \\ x^{k*} \end{pmatrix}$, где 1 — вектор-столбец, состоящий из единиц, а $x^* = (x_{ij})$, $i=1, \dots, m$; $j=1, \dots, n$, $x^{k*} = (x_{ij}^k)$, $i=1, \dots, m$; $j=1, \dots, n_k$.

Проверяемая гипотеза в таком случае имеет вид

$$H_0^* : \beta^{1*} = \beta^{2*} = \dots = \beta^{K*}.$$

Для проверки этой гипотезы придется снова вычислить остаток прогноза при справедливости нулевой гипотезы.

Так как

$$R_1^{k*} = (Y^k - b_0^k - b^{*k*})(Y^k - b_0^k - b^{*k*})',$$

то

$$R_1^* = \sum_{k=1}^K R_1^{k*} = \sum_{k=1}^K (Y^k - b_0^k - b^{*k*})(Y^k - b_0^k - b^{*k*})'.$$

При проверке гипотезы H_0^* при помощи F - или Λ -распределений следует число m заменить числом $m^* = m - 1$.

6. Вторая модификация поставленной задачи состоит в сравнении подмоделей.

Пусть I — некоторое подмножество множества $(0, 1, \dots, m)$, $I = (i_1, i_2, \dots, i_s)$. Обозначаем через $x(I)$ подвектор

$(x_{i_1}, \dots, x_{i_s})'$ вектора X и через $X(I)$ подматрицу $(x_{i_1} \dots x_{i_s})'$ матрицы X , где $x_{i_h} = (x_{i_h 1}, \dots, x_{i_h n})$ — i_h -ая строка матрицы X . Аналогично определяются и матрицы $X^k(I) = (x_{i_h, j}^k)$ ($h=1, \dots, s; j=1, \dots, n_k, k=1, \dots, K$). Также вводится подматрица $\beta(I)$ матрицы β , $\beta(I) = (\beta_{i_1, j_h})$ ($i=1, \dots, p; h=1, \dots, s$).

Подмоделью называется соотношение

$$y \approx E(y) = \beta(I)X(I),$$

в котором участвуют только аргумент-признаки x_{i_1}, \dots, x_{i_s} .

Подмодель в k -той генеральной совокупности имеет вид:

$$y^k \approx E(y^k) = \beta^k(I)X^k(I).$$

Задача сравнения подматриц состоит в проверке нулевой гипотезы

$$H_0(I) : \beta^1(I) = \beta^2(I) = \dots = \beta^K(I).$$

В случае $s=1$ данная методика позволяет проверять значимость каждого отдельного аргумент-признака в модели. Пошаговая реализация методики сравнения подмоделей позволяет выяснить оптимальный комплект аргумент-признаков, вызывающий различие моделей в разных подмножествах (в разные моменты времени).

При проверке гипотез следует отметить, что вместо $m+1$ придется применить число разных параметров s .

7. Если через J обозначается подмножество множества индексов всех исследуемых совокупностей, $J = (j_1, \dots, j_t)$, $J \subset (1, \dots, K)$, то простой модификацией описанной методики возможно проверить нулевую гипотезу

$$H_0(J) : \beta^{j_1} = \beta^{j_2} = \dots = \beta^{j_t}.$$

При вычислении параметров F - и Λ -распределений следует к

заменить на t .

Пошаговое применение этой методики полезно для выяснения самых сильно отличающихся друг от друга и от других подсовокупностей с точки зрения исследуемой модели.

8. В задачах анализа данных нередко целесообразно рассмотреть в качестве исходных данных не матрицы наблюдений X^k и X , а вычисленные по ним статистики – результаты первичной обработки. Такая ситуация имеет место, например, при весьма больших массивах данных, при подмассивах, исследуемых в разные времена или в разных местах.

Типичными результатами первичной обработки являются

- а) вектор–столбец оценок средних значений исследуемых признаков в k -той выборке \bar{z}^k , $\bar{z}^k = (\bar{x}_1^k, \dots, \bar{x}_m^k, \bar{y}_1^k, \dots, \bar{y}_p^k)'$;
- б) выборочная ковариационная матрица S^k признак–вектора z в k -той выборке,

$$S^k = \begin{pmatrix} S_{xx}^k & | & S_{xy}^k \\ \hline S_{yx}^k & | & S_{yy}^k \end{pmatrix};$$

- в) выборочная корреляционная матрица R^k признак–вектора z в k -той выборке,

$$R^k = \begin{pmatrix} R_{xx}^k & | & R_{xy}^k \\ \hline R_{yx}^k & | & R_{yy}^k \end{pmatrix}.$$

Соответствующие оценки \bar{z} , S^k и R^k по объединенной выборке, как правило, не известны.

Их легко найти по оценкам в отдельных выборках:

$$\bar{z} = \frac{1}{n} \sum_{k=1}^K n_k \bar{z}^k,$$

$$S = \frac{1}{n} \sum_{k=1}^K n_k s^k + \varepsilon ,$$

где

$$\varepsilon = \frac{1}{n} \sum_{k=1}^{K-1} \sum_{s=k+1}^K n_k n_s (\bar{z}^k - \bar{z}^s)(\bar{z}^k - \bar{z}^s)' ,$$

и

$$R = (\text{diag } S)^{-1/2} S (\text{diag } S)^{-1/2}$$

(здесь $\text{diag } A$ обозначает диагональную матрицу, в которой на диагонали находятся диагональные элементы матрицы A , а D^a при диагональной D есть матрица, в которой на диагонали элементы d_1^a , где d_1 — соответствующие элементы матрицы D).

9. При помощи выборочных ковариационных или корреляционных матриц остатки регрессии (4) и (5) вычисляются следующим образом.

В k -той подсовokuности имеем:

$$b^k = s_{yx}^k (s_x^k)^{-1}$$

и

$$R_0^k = s_y^k - s_{yx}^k (s_x^k)^{-1} s_{xy}^k ,$$

откуда получается суммарный остаток прогноза R_0 :

$$R_0 = \sum_{k=1}^K [s_y^k - s_{yx}^k (s_x^k)^{-1} s_{xy}^k] . \quad (6)$$

При справедливости нулевой гипотезы получим R_1 :

$$R_1 = s_y - s_{yx} (s_x)^{-1} s_{xy} . \quad (7)$$

Аналогичный результат можно получить и при помощи корреляционных матриц:

$$\begin{cases} \tilde{R}_0 = \sum_{k=1}^K [R_y^k - R_{yx}^k (R_x^k)^{-1} R_{xy}^k], \\ \tilde{R}_1 = R_y - R_{yx} (R_x)^{-1} R_{xy}. \end{cases} \quad (8)$$

Хотя остатки прогноза \tilde{R}_0 и \tilde{R}_1 , вычисленные по формулам (8) не совпадают с величинами R_0 и R_1 , найденными по формулам (6) и (7), отношение детерминантов, используемое при проверке гипотез, сохраняется.

10. Применение выборочных ковариационных или корреляционных матриц при вычислении регрессии полезно при неполных данных. Тогда все коэффициенты ковариации s_{ij}^k (корреляции) оценивают по всем наблюдениям из k -той совокупности, при которых i -ый и j -ый признак измерены. В таком случае получается оценка, которая максимально использует информацию, содержащуюся в выборке. При случайном расположении пробелов полученная оценка s^k является несмещенной, и полученные по ней оценки коэффициентов регрессии вполне применимы при анализе данных [2, 3].

Л и т е р а т у р а

1. Рао С.Р., Линейные статистические методы и их применения. Москва, 1968.
2. Тийт Э., Обработка неполных данных. Труды ВЦ ТТУ, 1977, 40, 14-20.
3. Koskel, S., Tiit, E., Regressioanalüüs. Programme kõigile XIV, 1978, 43-49.

КЛАССИФИЦИРОВАНИЕ НАУЧНЫХ ДАННЫХ С ПОЛУЧЕНИЕМ ЧАСТИЧНО ПОКРЫВАЕМЫХ КЛАССОВ

Р. Эзремаа

В настоящей статье мы затрагиваем проблемы классифицирования эмпирического материала, накопленного в ходе научного исследования, как биолого-медицинского, социолого-психологического, экономического, педагогического и т.п. характера.

Целью классифицирования научных данных является раскрытие сущности и закономерностей исследуемого материала с помощью представления массив информации в сжатом виде так, чтобы потеря информации не была чрезмерной.

Характерными чертами классифицирования научного материала являются: отсутствие обучающих выборок и априорной информации о характере распределения исследуемых признаков внутри классов и отсутствие жестких требований на результат классифицирования; но существование некоторых интуитивных представлений о результате (таких, как: разумное число классов и принадлежность некоторых объектов в один и тот же класс) — такая интуиция базируется на теоретических и на профессиональных соображениях.

Отметим, что вышеописанную задачу классификации можно назвать и задачей таксономии, а получаемые классы таксонами.

В связи с тем, что цель классифицирования здесь недостаточно четко определена, формализовать такую задачу нелегкая проблема.

1. Общие требования на постановку и решение задачи автоматической классификации научного материала

1.1. Разделение задач автоматической классификации

На основе работы [2] все множество задач автоматической классификации можно разбить на несколько основных групп:

1) по их постановке – на детерминистские и стохастические (или вероятностные) задачи;

2) по способу решения – на иерархические и неиерархические;

3) по наличию или отсутствию априорной информации о числе классов;

4) по тому, основывается ли решающий алгоритм непосредственно на состоянии объектов (т.е. на измерениях обследуемых признаков на каждом объекте) или на понятии меры близости (т.е. меры различия или сходства) между объектами.

Вышеприведенное разделение мы берем за основу при конкретизации постановки и решения задачи автоматической классификации научного материала.

1.2. Постановка задачи классифицирования научного материала

При задаче классифицирования научных данных обычно отсутствует такая достоверная априорная информация, которая нуждается в стохастическом подходе. Подходящим является здесь детерминистский подход, при котором исследуемый материал не

рассматривают как выборку из некоторой гипотической совокупности, а как самое совокупность. При таком подходе не существует эвристического этапа выбора гипотезы и "догадки" о закономерностях генеральной совокупности по конечной выборке.

О виде искомым таксонов надо сказать следующее. В ряде экспериментов было обнаружено, что при решении задачи таксономии человек использует простейшие (линейные) решающие функции. Поэтому, в работе [3] делается вывод, что если таксоны предназначены для непосредственного восприятия человеком, они должны иметь простую форму в виде гиперсферы.

Однако, то обстоятельство, что человек сам может осуществлять классифицирование в пространстве большой размерности только с помощью простейших решающих функций, не означает, что он не способен принять и понять более сложную классификацию.

Обычно эмпирический материал, накопленный в ходе научного исследования, обладает настолько "плохой" структурой, что неоправдан поиск классов простой формы с четко определенными классами. Хотя результат классифицирования научного материала должен быть упрощенным изложением материала, он кроме того должен довольно точно характеризовать и его внутреннюю структуру. Если на таком материале использовать методы нахождения непересекаемых классов простой формы, исследователь замечает при интерпретации получаемого результата, что расчленение некоторых объектов на разные таксоны с одной стороны является в соответствии с его теоретическим представлением о исследуемом материале, а с другой — эти таксоны заключают в себе "сомнительные" элементы, которые своим наличием

не поддерживают его гипотезу. Такое противоречие не позволяет ему ни принять новую гипотезу, ни опровергнуть старую.

При нахождении частично покрываемых таксонов произвольной формы такие "сомнительные" объекты обычно являются покрываемыми элементами между таксонами (на основе существующей информации они принадлежат к нескольким таксонам).

Классифицирование с получением пересекаемых классов соответствует и процессу мышления. Человек не всегда способен делать твердые, неопровержимые решения. На некотором уровне принятия решения он не в состоянии предпочесть одно решение другому, а при добавлении информации он иногда может такое предпочтение легко сделать. Отметим, что при классифицировании с получением частично покрываемых классов, добавочной информацией для решения вопроса о принадлежности "сомнительных" объектов к какому-то классу или нет, может быть и некоторая профессиональная информация исследователя (передать которую ЭВМ трудно).

Итак, задачу классифицирования научных данных поставим как детерминистскую задачу для нахождения классов произвольной формы, которые могут частично покрываться.

1.3. Решение задачи классифицирования научного материала

Эмпирический материал научных исследований по существу таков, что в нем разумно искать довольно однородные подмножества элементов, которые сравнительно изолированы от членов других подмножеств, и которые в то же время в свою очередь развиваются на сравнительно однородные, довольно изолированные друг от друга, подмножества.

Для того, чтобы получить точный и полный анализ исследуемого материала, с точки зрения учета его внутренней структуры, надо использовать иерархические процедуры. Иерархическое классифицирование дает исследователю возможность наглядной интерпретации проведенного анализа, но также и возможность выбирать уровень и соответствующее разбиение объектов в иерархическом дереве, исходя из своих теоретических соображений и из многих показателей, характеризующих разбиение.

При осуществлении иерархического классифицирования, которое основано на вышехарактеризованном представлении о классе как о множестве, не имеющего строго определенных границ, трудно представить иерархическое дерево традиционным образом. Поэтому, в таком случае результатом иерархического классифицирования представляются по очереди разбиения объектов на разных уровнях. Есть возможность выбирать на ЭВМ из таких уровней только один уровень (или несколько уровней) классификации соответственно каким-то критерием качества разбиения, в том числе соответственно наперед заданным числом классов. При частичной покрываемости между классами последняя значит задавание числа таких классов, в составе которых относительно мало покрываемых элементов.

1.4. Априорная информация о числе классов при задаче классифицирования научного материала

Ясно, что если мы реализуем классифицирование научного материала иерархическим процессом, то мы вместе с тем позволяем и задавание желаемого числа классов. Но в данном случае число классов является только одним из показателей, характе-

ризующих разбиение на каком-то уровне, а не представляет собой твердое условие на результат классифицирования, - т.е. процедур классифицирования не осуществляют, исходя из этого требования. Насколько произвольным может быть результат, когда в классифицировании исходят из наперед заданного числа классов, видно из примера, приведенного в следующей главе этой статьи.

1.5. Мера близости между объектами при задаче классифицирования научного материала

Мы еще не установили, должен ли решающий алгоритм основываться непосредственно на измерениях обследуемых признаков на объектах или на мере близости между объектами. На базе работ, в которых занимаются классифицированием первого вида, можно сказать, что обычно такое классифицирование осуществляется информационно-теоретическими средствами. Такая реализация сопровождается громоздкостью соответствующих процедур, большим удельным весом эвристических расчетов и необходимостью в емкой памяти ЭВМ.

При задаче классифицирования научного материала за основу берем подход, согласно которому описание любого объекта измерением m признаков интерпретируется геометрически в виде точки в m -мерном пространстве признаков. Естественно предположить, что геометрическая близость точек в этом пространстве означает сходство физических состояний соответствующих объектов. Поэтому, в данном подходе важное место приобретает понятие и мера близости объектов. Можно сказать, что, если имеется абстрактное множество, то введение понятия "близости"

между его элементами соответствует заданию в нем некоторой топологии. Тем самым пространство описаний превращается в топологическое пространство. Чтобы определить качество классификации, надо ввести функцию, которая указывала бы на степень согласованности данной классификации с введенной топологией.

Близость измеряют либо мерой различия, либо мерой сходства (о дефиниции таких мер например [4]). Надо отметить, что выбор такой меры на практике носит субъективный характер. Но как в работах [4] и [11] подчеркивается, при использовании одного и того же самого алгоритма классифицирования, самые разные меры различия (сходства) дают в практике почти одинаковые результаты. Но все-таки надо предпочитать такие методы классифицирования, которые позволяют наложение разной меры различия (сходства), определяющей структуру признакового пространства, методам, позволяющим наложение только функции расстояния, как это при выработанных методах часто бывает.

1.6. Процесс классифицирования научного материала и требования к методу классифицирования

Подытоживая теперь вышеприведенные обсуждения о классифицировании научного материала, можем сказать следующее. Процесс классифицирования будем рассматривать двухэтапно: к первому этапу относится конструирование матрицы различия (сходства) между классифицируемыми объектами; на втором этапе происходит конструирование классифицируемой системы, в которой объединяются объекты при различных уровнях различия (сходства), при том получаемые классы могут быть как непе-

ресекаемыми, так и частично покрываемыми. Такой двухэтапный процесс мы будем называть кластер-анализом.

Пусть у нас для конкретности данные для конструирования кластер-системы представлены в виде коэффициента различия (КР) над всеми исследуемыми объектами. Пусть их число n . Переход от КР к требуемому виду кластер-системы осуществляется применением кластер-метода к КР. При кластеризации научного материала разумно требовать от кластер-метода удовлетворения следующих условий:

- 1) независимость кластеризации от упорядочения объектов;
- 2) независимость кластеризации от шкалы;
- 3) если в данных некоторые кластеры четко вырисованы, то они должны сохраниться и после применения кластер-метода к исходному КР;
- 4) получаемый результат должен быть в некотором смысле самым лучшим при наложенных требованиях;
- 5) малым изменениям в данных должны соответствовать малые изменения в результатах.

Отметим, что в работах [9] и [12] представлена модель таксономии, в рамках которой можно математически корректно поставить и решить задачу классифицирования научного материала, так как все наши требования к классифицированию в нем удовлетворены.

2. Пример применения двух различных алгоритмов классифицирования, один из которых найдет классы, могущие частично покрываться

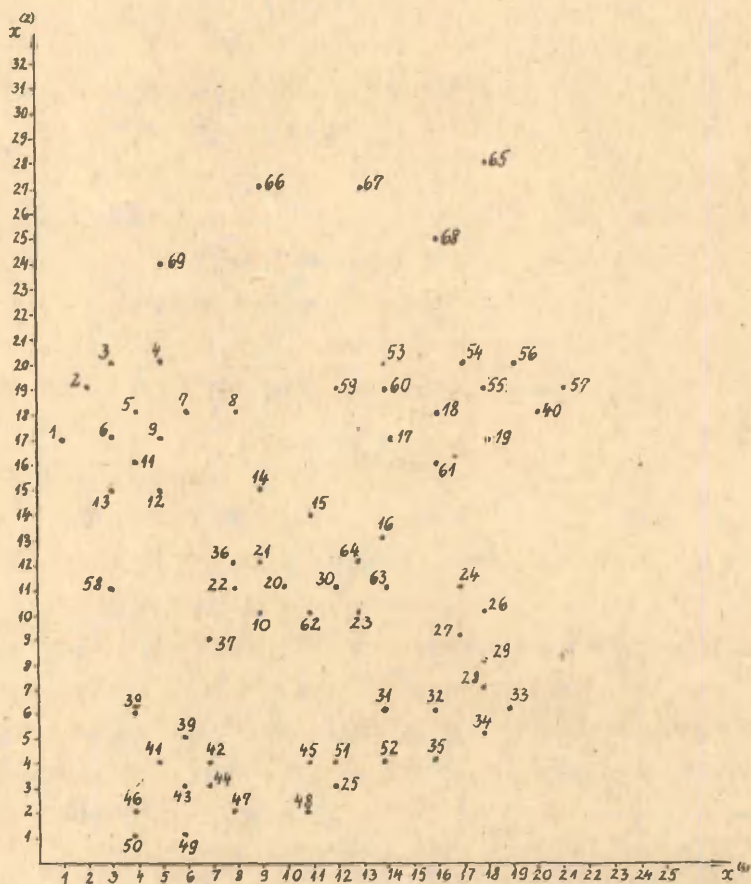
2.1. Исследуемая совокупность объектов и применяемые алгоритмы классифицирования

Пусть задана 69-элементная совокупность объектов, каждый из которых представлен в виде 2-мерного вектора измерений $X = (x^{(1)}, x^{(2)})$. В данном случае каждый объект можно считать точкой двумерного пространства. Соответствующие точки на плоскости могут быть расчленены человеком даже невооруженным глазом. Поэтому, на базе данной совокупности можно сравнивать как результаты применения разных методов классифицирования между собой, так и получаемые результаты с "истинным" расчленением объектов.

На рис. 1 исследуемые 69 объектов представлены в виде точек на плоскости. На этом рисунке визуально видно 5-6 скоплений наиболее близко расположенных друг от друга точек.

Демонстрируем, как выполняемость или невыполняемость некоторых вышехарактеризованных требований к классифицированию влияют на результат классифицирования. Мы сравниваем результаты при применении двух различных алгоритмов классифицирования, реализующих детерминистский подход с использованием меры близости между объектами. Первым является иерархический алгоритм ШК [6-7] для осуществления k -кластеризации, т.е. для получения кластеров, которые могут покрываться в объеме до $k-1$ объектов, принимая $1 \leq k \leq n-2$, где n - число исследуемых объектов. Вторым возьмем неиерархический алгоритм

Рис. 1. Представление 69 объектов точками на плоскости.



"ФОРЭЛЬ-2" [5]. В данном случае оба алгоритма используют мерой различия межобъектное эвклидовое расстояние. Отметим, что вообще алгоритм ППК может работать с любой мерой различия или сходства, в отличие от алгоритма "ФОРЭЛЬ-2", в котором исходными данными для построения классификации используют только межобъектные эвклидовы расстояния. При использовании алгоритма "ФОРЭЛЬ-2" надо задавать желаемое число классов, а при алгоритме ППК такая информация не нужна.

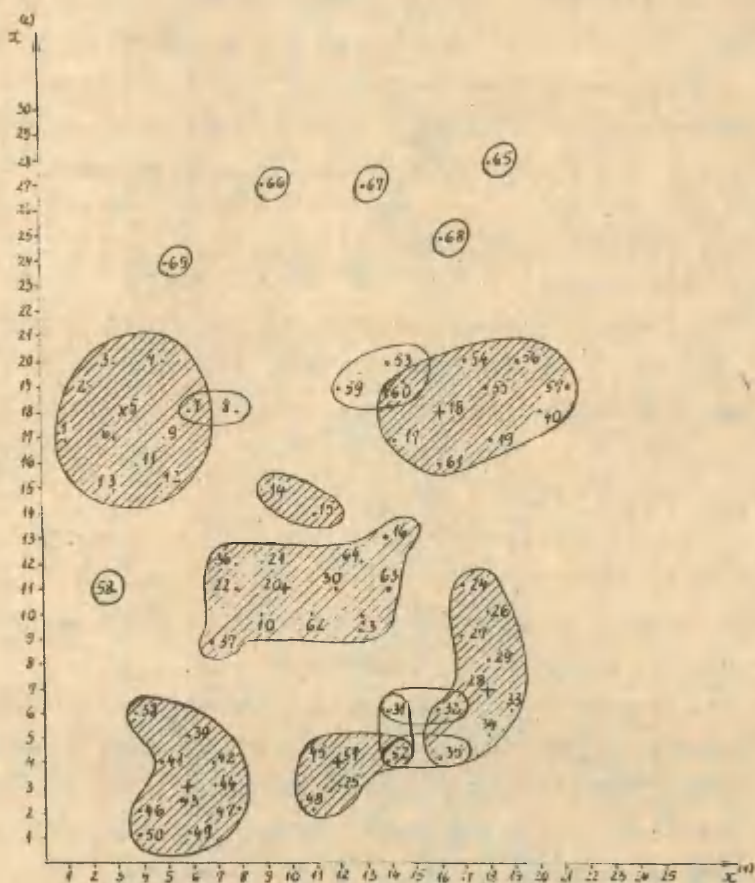
2.2. Применение алгоритма ППК

Применяя алгоритм ППК, можно получить кластер-систему, в которой объединяются объекты на нескольких уровнях различия (сходства). Есть возможность, что на ЭВМ выбирается из них т.н. уровень классификации, при котором получаемая кластер-структура в некотором смысле наиболее точно характеризует структуру данных.

На рис. 2 мы представляем результат 2-кластеризации алгоритмом ППК на уровне классификации, выбранном на ЭВМ. Отметим, что при 2-кластеризации получаемые кластеры могут частично покрываться до одного элемента.

В алгоритме ППК для каждого кластера вычисляется определенная мера внутренней связанности, которая может принимать значения в полуотрезке $(0, 1]$. При том значение 1 соответствует случаю, когда в полученном кластере ни один элемент не является покрываемым. На рис. 2 сильно-связанные кластеры (для которых в данном случае численные значения меры связанности превышают 0,833) представлены заштрихованными. Также на рис. 2 крестиками отмечены т.н. типичные представители

Рис. 2. Получаемые кластеры 2-кластеризации алгоритмом ППК на уровне классификации, выбранном на ЭВМ.



кластеров (число элементов которых больше двух). В качестве типичного представителя выбирается элемент, максимально связанный с остальными в данном кластере (на базе исходной матрицы различия или сходства). При применении алгоритма ШПК результатом получается 6 одноэлементных кластеров и 12 кластеров, число элементов которых больше одного.

Отметим, что алгоритм ШПК реализует метод классифицирования, который удовлетворяет всем требованиям, приведенным на стр. 62. Особое внимание обратим на то, что четковырисованные кластеры на рис. 1 сохранены после применения кластер-метода (см. рис. 2), и на то, что упорядочение объектов не изменит результат кластеризации.

2.3. Применение алгоритма "ФОРЭЛЬ-2"

Для сравнения результатов разных методов классификации применяем одну программу из хорошоизвестных алгоритмов группы "ФОРЭЛЬ", а именно "ФОРЭЛЬ-2". Алгоритм "ФОРЭЛЬ-2" производит расчленение совокупности объектов гиперсферами, — здесь реализуется характеризованная выше идея о конструировании таксонов простой формы. В информации к этой программе надо задавать желаемое число таксонов.

Кажется, что результат классифицирования, похожий на разбиение, найденное алгоритмом ШПК, можно получить при задании наперед числа (непересекаемых) таксонов 15. Исходя из такого задания, применение алгоритма "ФОРЭЛЬ-2" дает результатом разбиение на 15 таксонов, представленное на рис. 3. Для каждого выделенного таксона происходит поиск его типичной точки, т.е. точки, являющейся наиболее близкой к цен-

ТАКСОНОВ 15.



тру таксона. Эти типичные представители таксонов отмечены крестиками.

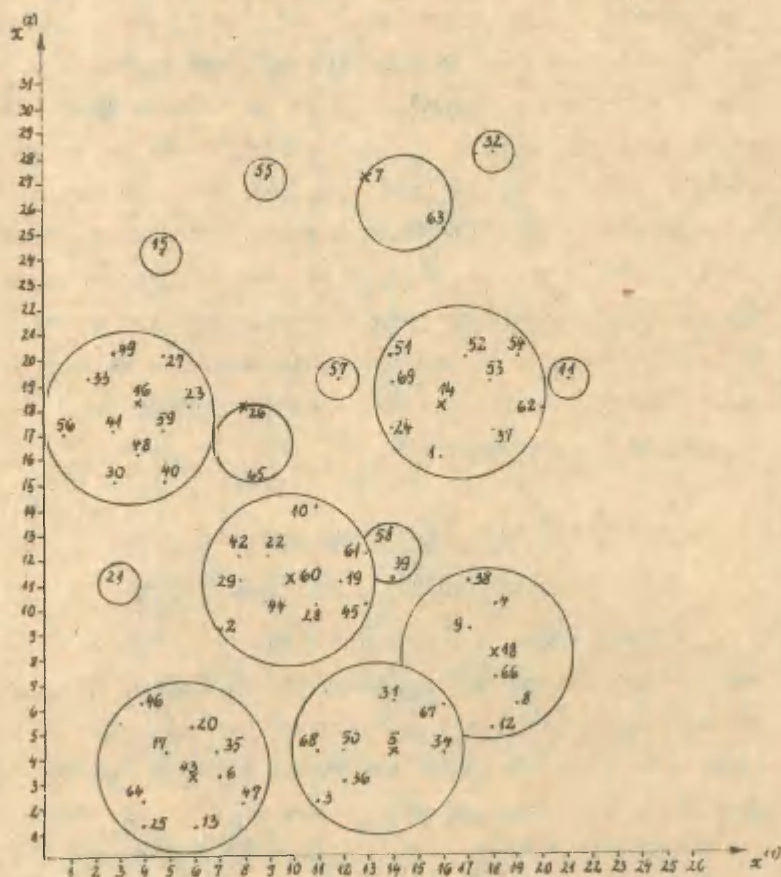
На рис. 3 видно, что таксоны, четко вырисованные на рис. 1, не сохранились после применения алгоритма "ФОРЭЛЬ-2". Это и естественно, потому что на метод классифицирования, который реализуется алгоритмом "ФОРЭЛЬ-2", не наложено требование о сохранении четко вырисованных таксонов в данных. Отметим, что на соответствующий метод не наложены ни требования, чтобы малым изменениям в данных соответствовали малые изменения в результатах, ни требования независимости результата от упорядочения объектов. Невыполняемость последнего условия иллюстрируем на рис. 4. На этом рисунке представлен результат классифицирования в случае изменения упорядоченности исследуемых объектов, и при неизменном желаемом числе таксонов 15. Видно заметное различие между классификациями, представленные на рисунках 3 и 4.

3. Подходящий метод классифицирования при конкретной задаче классифицирования

При применении алгоритма "ФОРЭЛЬ-2" мы увидели, что субъективное задание наперед числа искомых таксонов, принесет за собой и результат, субъективный в своей сущности. Поэтому алгоритм, в исходной информации которого задается желаемое число классов, не всегда дает хороших результатов в смысле наиболее точного отражения структуры исследуемой совокупности объектов в найденной классификации.

Вообще, при использовании любого метода классифицирова-

Рис. 4. Результат классифицирования в случае изменения упорядоченности исследуемых объектов, алгоритмом "ФОРЭЛЬ-2" при заданном числе таксонов 15.



ния, найденного в литературе, надо иметь в виду цель, на достижение которой соответствующее классифицирование направлено. С поставленной целью тесно связаны и требования, выполнение которых от метода классифицирования требуется. Так, например, мы требуем от метода классифицирования, который мы применяем на научные данные, выполнения ряда естественных условий (см. стр. 62). Но такие же требования не выполнены при любом публицированном методе классифицирования. Изложенный выше пример указывает, что принесет за собой, например, невыполнение условия независимости результата от упорядочения объектов.

Наш пример также показал, что при применении алгоритма ШПК получаемые классы (которые могли частично покрываться) довольно хорошо соответствовали нашим представлениям об истинной классификации. Отметим, что в литературе часто можно найти такие методы классифицирования, критерием качества которых служит только хорошая интерпретация найденных классов, согласованность полученной классификации с теоретическими представлениями при конкретной исследуемой совокупности данных. Довольно редко в литературе встречаем такой математически строгий подход к конструированию алгоритма классифицирования, который был основой для выработки алгоритма ШПК. Этот подход содержит следующие этапы:

- 1) точное математическое определение данных и результата классифицирования, позволяющее рассматривать метод классифицирования как преобразование из одной структуры в другую;

- 2) определение требований на такое преобразование;

3) определение существования, единственности и свойств метода, удовлетворяющего наложенным требованиям;

4) нахождение эффективного алгоритма.

Итак, нахождение пересекаемых классов при классифицировании научного материала оправдано и возможно, — что свидетельствует о существовании такого нетрадиционного способа классифицирования рядом с традиционным способом (нахождением непересекаемых классов). Подчеркиваем еще раз, что при классифицировании научного материала в качестве исходной информации не используются интуитивные представления исследователя об ожидаемом результате классифицирования (в том числе о желаемом числе классов). Мы реализуем математически строго представляемый метод классифицирования, который соответствует всем наложенным нами требованиям. Но так как мы притом требуем сохранения четко вырисованных кластеров, то получаемый результат, отражающий наиболее точно структуру совокупности объектов, в большинстве случаев совпадает и с интуитивными ожиданиями исследователя.

Л и т е р а т у р а

1. Дюран Б., Оделл П., Кластерный анализ. М., "Статистика", 1977.
2. Елисеева И.И., Рукавишников В.О., Группировка, корреляция, распознавание образов. М., "Статистика", 1977.
3. Загоруйко Н.Г., Методы распознавания и их применение. М., "Сов. Радио", 1972.

4. Ивахненко А.Г., Самообучающиеся системы распознавания и автоматического управления. Киев, "Техника", 1965.
5. Математическое обеспечение ЕС ЭВМ. Минск, 1976, 10, 6, 141-143.
6. Эзремаа Р., Общая теория конструирования кластер-систем и алгоритмы для нахождения их численных представлений. Труды ВЦ ТГУ, 1978, 42, 53-77.
7. Эзремаа Р., Алгоритм опознания кластеров к-дендрограммы. Труды ВЦ ТГУ, 1978, 42, 78-93.
8. Jardine, C.J., Jardine, N., Sibson, R., The structure and construction of taxonomic hierarchies. Math. Biosciences, 1967, 1, 173-179.
9. Jardine, N., Sibson, R., A model for taxonomy. Math. Biosciences, 1968, 2, 465-482.
10. Jardine, N., Sibson, R., The construction of hierarchic and non-hierarchic classifications. The Computer Journal, 1968, 11, 177-184.
11. Lerman, J.C., Les bases de la classification automatique. Paris Collection "Programmation", 1970.
12. Sibson, R., A model for taxonomy II. Math. Biosciences, 1970, 6, 405- 430.

ДВУХМЕРНЫЙ СТАТИСТИЧЕСКИЙ АНАЛИЗ ДАННЫХ

Ю. Вилисмья

Важной частью системы статистического анализа данных в ВЦ ТТУ является т.н. двухмерный статистический анализ.

Комплект соответствующих программ работает с 1976 года как часть системы для ЭВМ "Минск-32". Теперь аналогичная методика реализована и для ЭВМ ЕС-1022 для работы в операционной системе ДОС.

Ниже излагается краткая характеристика этой методики, причем предполагается знание основных принципов системы статистической обработки данных (см. [1]).

1. Общие положения

Пусть задана совокупность исследуемых признаков $Z = \{x_1, \dots, x_m\}$. Двухмерный статистический анализ этих признаков состоит в совместном исследовании каждого признака из одной группы признаков

$$Z_1 = \{x_{1_1}, \dots, x_{1_r}\} \\ (x_{1_s} \in Z; \quad 1 \leq 1_s \leq m; \quad s=1, \dots, r)$$

вместе с каждым признаком из второй группы

$$Z_2 = \{x_{k_1}, \dots, x_{k_p}\}$$

$$(x_{k_t} \in Z; \quad 1 \leq k_t \leq m; \quad t=1, \dots, p),$$

причем эти группы могут частично или полностью совпадать.

Основным элементом двумерного статистического анализа является анализ одной пары признаков ($X_1 \in Z_1; X_2 \in Z_2$). Это означает образование таблицы частот и вычисление эмпирического и условных распределений рассматриваемых признаков X_1 и X_2 , а также вычисление некоторых параметров этого распределения, в том числе показателей связи.

Двухмерный статистический анализ всех исследуемых признаков осуществляется в ЭВМ при помощи последовательности анализов отдельных пар признаков. В конце анализа коэффициенты связи при надобности соединяются в $(r \times p)$ -матрицу и печатаются.

Для простоты в дальнейшем будем обозначать $X_1 = X$ и $X_2 = Y$.

2. Совместное исследование пары признаков X и Y

Пусть заданы выборки значений (x_1, \dots, x_N) и (y_1, \dots, y_N) признаков X и Y соответственно. Предположим, что признак X имеет k классов распределения: A_1, \dots, A_k , а Y l классов: B_1, \dots, B_l ($k, l \geq 1$).

Классом распределения в принципе может быть любое связанное множество значений данного признака. Разделение области значений на классы задается при помощи данных классифицирования (см. [4], стр. 30-42). Классы выбираются непересекающимися друг друга, причем ими необязательно охватывать все

элементы выборки.

Рассмотрим теперь отдельные части листинга, полученного в результате анализа пары признаков X и Y .

2.1. Заголовок. Выпечатываются имена и номера признаков X и Y , уровень значимости и другие общие параметры обработки.

2.2. Таблица частот. В таблице частот печатается матрица $\{n_{ij} \mid i=1, \dots, k; j=1, \dots, l\}$, где n_{ij} - количество объектов (индивидов), попавших одновременно в класс A_i признака X и в класс B_j признака Y .

2.3. Таблица совместного эмпирического распределения признаков. В таблице печатаются относительные частоты (в процентах) $\left\{ \frac{n_{ij}}{n} 100 \right\} (i=1, \dots, k; j=1, \dots, l)$, где $n = \sum_{i=1}^k \sum_{j=1}^l n_{ij}$. (Очевидно, $n \leq N$.)

2.4. Параметры распределения. Вычисляются эмпирические средние:

$$\bar{x} = \frac{1}{n} \sum_1 x_1, \quad \bar{y} = \frac{1}{n} \sum_1 y_1$$

и эмпирические стандартные отклонения:

$$s_X = \sqrt{\frac{1}{n-1} \sum_1 (x_1 - \bar{x})^2}, \quad s_Y = \sqrt{\frac{1}{n-1} \sum_1 (y_1 - \bar{y})^2}$$

(в них и в формуле вычисления коэффициента линейной корреляции, суммирование производится по всем значениям индекса i при которых $x_1 \in \bigcup_{s=1}^k A_s$ и $y_1 \in \bigcup_{t=1}^l B_t$), а также доверительные границы эмпирических средних при заданном уровне значимости α :

$$\bar{x} - \frac{S_X t_\alpha}{\sqrt{n}}, \quad \bar{x} + \frac{S_X t_\alpha}{\sqrt{n}} \quad \text{и} \quad \bar{y} - \frac{S_Y t_\alpha}{\sqrt{n}}, \quad \bar{y} + \frac{S_Y t_\alpha}{\sqrt{n}},$$

где t_α — соответствующий квантиль t -распределения Стьюдента.

Кроме того вычисляются параметры связи:

коэффициент линейной корреляции:

$$r = \frac{\sum_1 (x_1 - \bar{x})(y_1 - \bar{y})}{(n-1) S_X S_Y},$$

χ^2 -статистика:

$$\chi^2 = n \left(\sum_{i=1}^k \sum_{j=1}^l \frac{n_{ij}^2}{n_{i.} \cdot n_{.j}} - 1 \right),$$

где $n_{i.} = \sum_{j=1}^l n_{ij}$ и $n_{.j} = \sum_{i=1}^k n_{ij}$,

коэффициент связи Чупрова:

$$T = \sqrt{\frac{\chi^2}{n \sqrt{(k^*-1)(l^*-1)}}},$$

где k^* (соответственно l^*) — число непустых строк (столбцов) в таблице частот,

обе корреляционные отношения:

$$\eta(Y/X) = \left[\sum_{i=1}^k \frac{n_{i.} (\bar{y}_i - \bar{y})^2}{S_Y^2 (n-1)} \right]^{1/2},$$

$$\eta(X/Y) = \left[\sum_{j=1}^l \frac{n_{.j} (\bar{x}_j - \bar{x})^2}{S_X^2 (n-1)} \right]^{1/2},$$

где \bar{y}_i ($i=1, \dots, k$) и \bar{x}_j ($j=1, \dots, l$) — эмпирические условные средние признаков Y и X соответственно. (Формулы для их вы-

числения задаются ниже.)

При каждом параметре связи устанавливается и его статистическая значимость.

Отметим, что все перечисленные величины вычисляются в данном комплекте программ непосредственно по измеренным значениям признаков. В комплекте программ для ЭВМ "Минск-32" при вычислении этих величин все наблюдения признака, находящиеся в i -ом классе, заменяются т.н. представителем этого класса, которым в общем случае является взвешенное среднее тех элементов выборки, которые попадают в данный класс.

2.5. Информация о соединении классов. Для вычисления корреляционных отношений и коэффициента корреляции Чупрова в таблице частот последовательно соединяются те классы, в которых мало индивидов (меньше заданного числа ω), со соседними классами, пока не достигается достаточная величина класса. На листинге показывается, какие классы соединялись.

2.6. Условные распределения признаков. Условные распределения признака X определены для каждого класса B_j ($j=1, \dots, l$) признака Y . Эти условные распределения расположены в таблице условных распределений по столбцам, причем номер столбца указывает номер класса B_j , который определяет условие. Таким образом, в j -ом столбце ($j=1, \dots, l$) располагаются величины $100 \frac{n_{ij}}{n_{.j}}$ ($i=1, \dots, k$). Очевидно, сумма величин в каждом столбце равняется единице (100%).

Аналогично вычисляются условные распределения признака Y при условиях $X \in A_i$ ($i=1, \dots, k$). Они располагаются в соответствующей таблице по строкам.

2.7. Параметры условных распределений. Параметрами условных распределений признака X являются эмпирические условные средние $\bar{x}_j = EX(Y \in B_j)$ (т.е. эмпирические средние признака X при условии, что значения признака Y относятся к классу B_j):

$$\bar{x}_j = \frac{1}{n_{\cdot j}} \sum_{y_1 \in B_j} x_1 \quad (j=1, \dots, l),$$

эмпирические условные стандартные отклонения:

$$s_{X(j)} = \sqrt{\frac{1}{n_{\cdot j} - 1} \sum_{y_1 \in B_j} (x_1 - \bar{x}_j)^2} \quad (j=1, \dots, l).$$

Вычисляются также доверительные границы условных средних признака X :

$$\bar{x}_j - \frac{s_{X(j)} t_{\alpha}}{\sqrt{n_{\cdot j}}} \text{ и } \bar{x}_j + \frac{s_{X(j)} t_{\alpha}}{\sqrt{n_{\cdot j}}} \quad (j=1, \dots, l).$$

Все эти параметры выпечатываются вместе с таблицей эмпирического распределения (см. 2.3) в виде дополнительных четырех строк под таблицей.

Аналогично определяются параметры условных распределений признака Y . Они выпечатываются в форме дополнительных столбцов к той же таблице распределения.

2.8. Регрессионные линии. Выпечатываются графики регрессионных линий (т.е. графики условных средних) признаков X и Y относительно признаков Y и X соответственно (см. [2]).

2.9. Корреляционное поле. Выпечатывается корреляционное поле признаков X и Y (см. [2]).

2.10. Номера индивидов. Выпечатываются номера тех индивидов, соответствующие значения которых попали в таблицу частот. Какой-то i -ый индивид может не попасть в таблицу, во-первых, потому что значение x_i не включается ни в одном из классов A_j или y_i ни в одном B_j ; во-вторых, потому что данные неполны и значение X или Y не измерено при i -ом индивиде.

Все перечисленные части листинга при помощи ключа режима выпечатки можно либо выпечатать, либо опускать независимо друг от друга (за исключением параметров условных распределений, которые печатаются всегда вместе с таблицей эмпирического распределения). Это дает заказчику возможность заказывать только интересующую его часть информации.

3. Заказ двумерного статистического анализа

Исходом двумерного статистического анализа является массив значений признаков X_1, \dots, X_m в стандартном виде на магнитном диске (см. [3]). Двумерный анализ происходит на основе заказа, в котором нужно указывать:

- 1) перечень исследуемых признаков,
- 2) данные классифицирования,
- 3) общие параметры для анализа,
- 4) режим выпечатки.

Заказ вводится в ЭВМ через перфокарточное устройство ввода.

Перечень исследуемых признаков задается в виде

$$i_1, \dots, i_r * k_1, \dots, k_p), \quad (1)$$

где i_s ($s=1, \dots, r$) – номера признаков из группы Z_1 , а k_t ($t=1, \dots, p$) – номера признаков из группы Z_2 . Если $Z_1 = Z_2$, то перечень задается в виде

$$i_1, \dots, i_r) . \quad (2)$$

Признаки могут быть любого типа: А, В, С или R (см. [3]).

Данные классифицирования задаются только для тех признаков, для которых надо изменить имеющиеся в стандартном массиве признаков данные классифицирования.

Заказчик должен задавать некоторые общие параметры работы:

уровень значимости α ,

имя массива значений признаков на магнитном диске,

количество индивидов N в этом массиве,

критерий пустого класса ω (см. 2.5),

критерий пустой таблицы ε .

Таблица частот считается пустой, если $n < \varepsilon$. В таком случае двумерный анализ рассматриваемой пары признаков не осуществляется.

Все результаты двумерного статистического анализа данных выводятся только на листе печатающего устройства. Листинг двумерного анализа признаков X_1, \dots, X_m состоит из последовательности листингов двумерных анализов отдельных пар признаков. Если перечень признаков задан в виде (1), то эта последовательность состоит из анализов всех пар признаков

$$(X_{i_s}, X_{k_t}) \quad \begin{aligned} &(X_{i_s} \in Z_1; s=1, \dots, r; \\ &X_{k_t} \in Z_2; t=1, \dots, p). \end{aligned}$$

Если перечень задан в виде (2), то эта последовательность состоит из анализов пар

$$(x_{1_s}, x_{1_t}) \quad (x_{1_s}, x_{1_t} \in Z_1; 1_s \geq 1_t) .$$

Притом объем выпечатываемой информации при каждой паре определяется режимом выпечатки.

Л и т е р а т у р а

1. Труды ВЦ ТГУ, 1977, 40.
2. Вилисмяз Ю., Выпечатка графиков первичного статистического анализа данных. Труды ВЦ ТГУ, 1978, 42, 117-124.
3. Коскель С., Структура данных на этапе ввода. Труды ВЦ ТГУ, 1978, 42, 105-111.
4. Programme kõigile XI, Tartu, 1976.

ЭМПИРИЧЕСКОЕ ИССЛЕДОВАНИЕ ГОРОДСКОЙ СРЕДЫ ПРИ ПОМОЩИ СИСТЕМЫ СТАТИСТИЧЕСКОЙ ОБРАБОТКИ ДАННЫХ В ВЦ ТТУ

Л.М. Тоодинг

В 1978–1979 году в вычислительном центре ТТУ проводилась совместная с горисполкомом г. Тарту научно-исследовательская работа по эмпирическому анализу городской среды. Целью исследования было составление математико-статистических моделей восприятия некоторых аспектов городской среды, а также анализ разрабатываемых моделей.

Обработка соответствующего эмпирического материала производилась на ЭВМ "Минск-32" и ЕС-1022, используя систему статистической обработки данных в ВЦ ТТУ (см. [6]) и пакет SSP (дополненная нами версия, см. [5]), притом практически в полном объеме. Поэтому, настоящему обзору одной конкретной прикладной работы, проведенной в ВЦ ТТУ, можно присвоить также определенное методологическое значение в смысле иллюстрации возможностей имеющегося у нас программного обеспечения.

1. Структура данных и выбор основных рабочих гипотез

1.1. Эмпирические данные рассматриваемого исследования собирались в виде анкет, заполнявшихся респондентами при по-

мощи интервьюеров. Респонденты выбирались случайно, притом попадание в выборку являлось равновероятным для всех тар-туанцев в возрасте от 18 до 75 лет. Из 1100 разданных анкет пригодными для анализа оказались 1068.

За заполнением анкет следовала их проверка и кодирование, а затем ввод в фонд данных системы статистического анализа. Из 160 вопросов анкеты путем формализации образовался массив данных с 640 признаками.

Напомним, что основной структурной единицей данных при системе является прямоугольная матрица $X = (x_{ij}, i=1, \dots, M; j=1, \dots, N)$, представляющая выборку объема N M -мерного признак-вектора $\bar{X}_M = (X_1, \dots, X_M)$, где элемент x_{ij} — значение i -го признака у j -го индивида (объекта). Данные могут быть неполными и иметь произвольный статистико-вероятностный характер (см. [7]).

В рассматриваемой задаче средняя доля пробелов образует 11% от объема выборки, а распределение признаков по типам следующее:

качественные — 19%

количественные упорядоченные — 66%

количественные — 15%.

Таким образом, в основном обработке можно было провести методами анализа количественной информации (осмысляя притом результаты упорядоченных признаков качественно), значит, данные соответствуют тому программному обеспечению, которым мы пользовались.

1.2. Содержательные гипотезы данной работы затрагивают в основном пространственно-деятельные соотношения в городской среде. В этих целях были выделены определенные пространственные зоны и типы действия. Рассматривалось следующее разделение городской среды:

- 1) квартира, жилое помещение;
- 2) жилище, жилой дом;
- 3) окружность дома;
- 4) соседство;
- 5) часть города, микрорайон;
- 6) город в целом;
- 7) загород.

В причинном плане, для установления соотношений детерминации, данное разделение осмысливается как некоторая иерархия и используется формальной базой для выведения основных рабочих гипотез. Притом есть смысл рассматривать оба направления – сверху вниз и снизу вверх, – считая все предыдущие ступени аргументами для данного.

Так, например, стереотип микрорайона мыслимо объяснить как характером связей с загородом и с городом в целом, так и непосредственно квартирой, отношениями к первичной жилой среде.

1.3. Большинство признаков рассматриваемого массива данных и служит для описания пространственных зон, но наряду с этим рассматриваются также следующего рода действия:

- 1) движение по маршруту "дом – место работы";
- 2) действия в домашнем хозяйстве;

- 3) движение в целях обслуживания, закупок;
- 4) занятия физического отдыха;
- 5) занятия отдыха нефизического характера.

Основные рабочие гипотезы при данном разделении выражаются в зависимости 2-5 от 1 и во взаимосвязях 4 и 5, остальные возможные соотношения в признаках не отражаются.

1.4. Действия естественным образом прикрепляются к определенным зонам. В данной работе рассматривалось следующее их совместное распределение.

	Работа	Дома	Обслу- живание	Физи- ческий отдых	Нефизи- ческий отдых
Квартира					
Дом					
Окружность					
Соседство					
Часть города					
Город в целом					
Загород					

Приведенной схемой выражаются и выдвинутые и проверенные в данном исследовании рабочие гипотезы о связи (как правило, о взаимосвязи) между деятельным стереотипом и пространственными воображениями.

Фоном всех действий в соответствующих зонах являются социально-демографическая характеристика респондента и семьи, стереотип движения по городу и снабженность городской информацией.

1.5. Отметим еще один способ содержательного классифицирования данного материала. Именно, все названные выше блоки по существу распадаются на две части: на описание нынешнего положения ("существующее") и на описание идеального положения ("желаемое"). Различием в смысле этих двух аспектов характеризуется по соответствующим признакам степень удовлетворенности при каждой пространственной зоне и при каждом действии.

2. Методика статистического анализа

2.1. Эмпирический анализ городской среды методами математической статистики характеризуется довольно скромным практическим опытом, к тому же прибавляются и существенные локальные особенности. Проведенные до сих пор в г. Тарту эмпирические исследования (напр. [2,3]) хотя и пошли впрок своей ценной информацией об отдельных сторонах городского организма, но не затрагивали городской системы в целом, как это ожидается при рассматриваемой работе.

В методологическом плане такое отсутствие априорных знаний означает, с одной стороны, сравнительно формальную, по некоторому принципу систематическую постановку рабочих гипотез (как это и указывалось в 1), а с другой — намерение воспользоваться абсолютно всей мощностью имеющихся в нашем распоряжении средств обработки данных и соблюдения принципа максимального сохранения исходной информации на всех этапах анализа.

2.2. На этапе ввода в целях наименьшей потери информации проводилась сравнительно слабая формализация данных. Производилось лишь закодирование в численный вид, окончательная же формализация была отложена на этап обработки. Как правило, в течение анализа для одной и той же группы признаков использовалось несколько разных схем формализации.

Уточняем, что под формализацией в данном случае имеются в виду разнообразные преобразования признаков: перекодирование, вычисления по правилу типа арифметического выражения, индексирование и т.п. операции в системе статистической обработки данных (см. [4], стр. 109).

Например, на основе информации, представленной респондентом о структуре желаемой им квартиры (величина и функция каждого помещения), выводится функциональный тип квартиры. Притом количество названных помещений может быть произвольным.

В рассматриваемом массиве данных преобразованиям в целях формализации подвергалось 45% исходных признаков, притом в непреобразованном виде с ними был проведен в основном лишь первичный анализ.

2.3. В силу слабой формализованности и сложности изучаемого (объекта), а также с преднамерением затронуть "кое-что обо всем", в данном случае рассматривается пространство признаков сравнительно высокой размерности. Поэтому, общую схему анализа можно характеризовать как обработку исходных признаков отдельными блоками, притом связь между блоками рассматривается по сводным признакам блока. Точнее это значит:

следующее.

Пусть задано некоторое разбиение множества элементов признак-вектора \bar{X} на P непересекающихся множеств, которым соответствуют признак-векторы $\bar{X}_{M_1}^1, \dots, \bar{X}_{M_P}^P$, где $\sum_{i=1}^P M_i = M$. Этим разбиением заодно определяется и P подматриц от выборки X , являющихся исходом для операций системы.

Разбиение может определяться, например, содержательной структурой признаков. В данном случае и рассматриваются совокупности признаков, соответствующие перечисленным в 1 зонам и действиям, к чему прибавляются еще блоки социально-демографических данных респондента и его семьи, блок о городской информации и способах движения.

В пределах каждого блока проводилась обработка с использованием, как правило, всех операций системы (см. [1]). Составлялись эмпирические распределения (одного признака, пары признаков, признак-вектора произвольной размерности), проводилось множественное сравнение средних методом Шеффе, осуществлялся анализ корреляционных матриц, применялись множественные дисперсионный и регрессионный анализы.

В целях сжатия информации производились разные преобразования признаков и применялся факторный анализ. В результате этого i -ый блок $\bar{X}_{M_i}^i$ для дальнейшего анализа представлялся сводным признак-вектором $\bar{Y}_{m_i}^i$, где $m_i \ll M_i$. Анализом признак-вектора $\bar{Y} = (\bar{Y}^1, \dots, \bar{Y}^P)$ осуществлялась проверка гипотез о взаимосвязях разных содержательных блоков, притом использовалась та же методика, что и внутри блоков.

Например, составлялись модели зависимости факторов квартиры от факторов способа проведения досуга и факторов соот-

ношений с загородом. Факторы желаемого типа района города прогнозировались через факторы семьи респондента и факторами жилища и т.д.

Таким подходом удалось в множественный анализ признаков включить косвенно и такие объемистые совокупности исходных признаков, совместный анализ которых был бы технически неосуществим. На уровне сводных признаков наблюдалась более высокая множественная связь в сравнении с любыми надлежащими моделями исходных признаков. Вдобавок облегчается в некоторых случаях (например, в регрессионных моделях) интерпретация результатов, так как получающиеся факторным анализом сводные признаки взаимоортогональны и отпадает надобность в осмыслении совлияний.

Переход на уровень сводных признаков в свою очередь можно осуществить и для некоторого разбиения \bar{Y} и т.д.

2.4. Наряду с естественной внутренней структурой совокупности признаков в данном исследовании рассматривались и некоторые способы априорной классификации индивидов. При этом имеется в виду не только анализ подматриц от X , соответствующих индивидам, при которых выполнено заданное логическое условие (см. [4], стр. 107), но и анализ данных об определенных сводных объектах. Хотя непосредственно измерялся респондент, житель города, и в отношении к нему и была устроена соответствующая статистическая выборка, согласно целям исследования, потребовалось также рассмотрение семьи и части города как самостоятельных объектов.

Данные для этих сводных объектов сгенерировались от вы-

борки X по информации соответствующих рассматриваемому суперобъекту респондентов.

Признаки семьи вычислялись операциями преобразования признаков: усреднением, индексированием, кодированием по логическому условию (см. [4]). Объекты "респондент" и "семья" находятся в однозначном соответствии, так как по правилу составления выборки весьма маловероятно попадание в нее нескольких представителей одной и той же семьи. Поэтому, получающиеся сводные данные о семье включались в первоначальный массив и обрабатывались вместе со всеми исходными признаками.

В качестве данных семьи вычислялись, например, демографический и социальный тип семьи, образовательный статус, средние показатели дохода, возрастной индекс и т.п.

Значения признаков для объекта "часть города" вычислялись операцией усреднения по значениям соответствующих признаков у респондентов рассматриваемой части города. Массив данных объекта "часть города" не совместим с массивами респондента и семьи, поэтому обрабатывался отдельно. В основном применялись методы классификации объектов. Так как передача информации в данном случае (в условиях нашего программного обеспечения) осуществляется повторным вводом, а не системно, то суперобъектов типа "часть города" следует рассматривать лишь в порядке исключения.

Вышесказанным охватываются основные характерные черты данной задачи. Опыт применения системы статистической обработки данных для решения ее убедил в том, что методологичес-

кие основы системы выбраны удачно и позволяют путем сравнительно формального подхода даже при скромной априорной информации получить содержательно хорошо осмысляемые результаты.

Л и т е р а т у р а

1. Прээм М., Тоодинг Л.М., Операции и принципы их следования. Труды ВЦ ТТУ, 1977, 40, 115-135.
2. Социальные предпосылки формирования пространственно-предметной среды. Тарту, 1972.
3. Типология уклада жизни городских и сельских семей на основе социологических и пространственно-предметных характеристик. Тарту, 1974.
4. Тоодинг Л.М., Преобразование стандартной матрицы. Труды ВЦ ТТУ, 1977, 40, 106-114.
5. Тоодинг Л.М., Применение пакета SSP в системе STR. Труды ВЦ ТТУ, 1978, 42, 147-156.
6. Тоодинг Л.М., Система статистической обработки данных в вычислительном центре ТТУ. Труды ВЦ ТТУ, 1977, 40, 3-7.
7. Тоодинг Л.М., Структура данных. Труды ВЦ ТТУ, 1977, 40, 75-85.

О ПРОЕКЦИОННЫХ МОДЕЛЯХ МНОГОМЕРНОГО ИНВАРИАНТНОГО ШКАЛИРОВАНИЯ

А.Б. Хмельницкая

При решении различных вопросов, связанных с обработкой экспериментальных данных, описываемых большим числом наблюдаемых характеристик, возникает проблема "сжатия" имеющейся эмпирической информации, т.е. задача адекватного описания изучаемых объектов посредством существенно меньшего набора интегральных признаков. Для решения подобного рода задач могут быть применены методы инвариантного шкалирования [1-4].

Наличие различных моделей инвариантного шкалирования позволяет изучать эмпирические данные различными способами и тем самым получать о них более объективную сжатую информацию.

В настоящей работе исследуются проекционные модели многомерного инвариантного шкалирования. В [4] показано, что все проекционные многомерные модели одной размерности при фиксированной функции близости в исходном пространстве признаков в случае, когда в качестве функции близости на выстраиваемой многомерной шкале рассматривается квадрат евклидовой метрики, эквивалентны и в следствии теоремы существования для соответствующей метрической модели [2] имеют решение (существуют функции шкалирования). В силу эквивалентности раз-

личных моделей достаточно рассмотреть только одну многомерную проекционную модель. Оказывается, что основные факты, справедливые для проекционных одномерных моделей [3-4] и для многомерной метрической модели [2], остаются в силе. А именно:

1) искать решение поставленной задачи можно двумя различными способами, что в связи с применением в вычислительных процедурах локальных методов дает возможность использовать решение, найденное одним путем, в качестве начального приближения для получения уточненного решения с помощью другой процедуры;

2) в случае вырожденности функции близости в исходном пространстве признаков вектор-функция шкалирования допускает два различных конечнопараметрических представления, причем число параметров обоих представлений не зависит от объема выборки, а лишь от "сложности" задачи;

3) для отыскания параметров, определяющих вектор-функцию шкалирования, для обоих представлений получены системы нелинейных уравнений;

4) в ряде случаев, в частности, в случае статистической независимости строящихся шкал, даются аналитические представления вектор-функции проекционного шкалирования.

Пусть (X, μ, r) - исходная генеральная совокупность объектов. Полагаем $r \in L^\infty(X \times X, \mu \times \mu)$. r - функция близости в R^n :

$$r(x, y) = \sum_{i=1}^n (x_i - y_i)^2 \quad \text{для любых } x, y \in R^n. \quad (1)$$

Для любого отображения $f: X \rightarrow R^n$ определим

$$\varphi_f(z, w) = \varphi(f(z), f(w)) \quad \text{для любых } z, w \in X. \quad (2)$$

Через (\cdot, \cdot) и $\|\cdot\|$ обозначим соответственно скалярное произведение и норму в $L^2(X \times X, \mu \times \mu)$.

Так как φ_f , определяемая формулами (1) – (2), положительно однородна степени 2 ([4]), то при исследовании проекционных моделей n -мерного инвариантного шкалирования достаточно рассмотреть задачу построения отображения $t: X \rightarrow R^n$ класса $L^4(X, \mu; R^n)$, доставляющего глобальный максимум функционалу

$$\psi_n^{(n)}(f) = \frac{(f, \varphi_f)}{\|\varphi_f\|} = \frac{\int_X \int_X f(z, w) \sum_{i=1}^n [f_i(z) - f_i(w)]^2 d\mu(z) d\mu(w)}{\sqrt{\int_X \int_X \left\{ \sum_{i=1}^n [f_i(z) - f_i(w)]^2 \right\}^2 d\mu(z) d\mu(w)}} \quad (3)$$

на множестве отображений $f: X \rightarrow R^n$ класса $L^4(X, \mu; R^n)$, удовлетворяющих условию

$$\|\varphi_f\|^2 = \int_X \int_X \left\{ \sum_{i=1}^n [f_i(z) - f_i(w)]^2 \right\}^2 d\mu(z) d\mu(w) = \|f\|^2. \quad (4)$$

Заметим, что функционал $\psi_n^{(n)}(f)$ и условие (4) инвариантны относительно сдвига f на произвольный постоянный вектор. Поэтому исследуемую задачу проекционного многомерного инвариантного шкалирования достаточно рассматривать на подпространстве $L^{(n)} \subset L^4(X, \mu; R^n)$

$$L^{(n)} = \{f \in L^4(X, \mu; R^n) \mid \bar{f} = 0\},$$

$$\text{где } \bar{f} = (\bar{f}_1, \dots, \bar{f}_n); \quad \bar{f}_i = \int_X f_i d\mu \quad 1 \leq i \leq n.$$

Кроме того, аналогично метрической модели многомерного инвариантного шкалирования, функционал $\psi_n^{(n)}(f)$ и условие (4) ин-

вариантны относительно вращения вектора f в R^n .

Как показано в [4], поставленную задачу можно решать двумя способами.

1. Т.к. для любого $a \neq 0$ $\Psi_{\Pi}^{(n)}(af) = \Psi_{\Pi}^{(n)}(f)$, то ищем любую вектор-функцию $t(z)$, на которой достигается глобальный максимум функционала $\Psi_{\Pi}^{(n)}(f)$ на множестве $f \in L^{(n)}$, а затем строим $t'(z) = at(z)$, причем величину a выбираем так, чтобы вектор-функция $t'(z)$ удовлетворяла условию (4).

Эту задачу будем кратко называть задачей $\Pi^{(n)}_1$.

2. Решаем задачу на условный максимум функционала

$$\Psi_{\Pi}^{(n)}(f) = (r, f_f) = \int_X \int_X r(z, w) \sum_{i=1}^n [f_i(z) - f_i(w)]^2 d\mu(z) d\mu(w)$$

при условии (4).

Эту задачу будем называть задачей $\Pi^{(n)}_2$.

Хотя так же как и в одномерном случае решения задач $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$ совпадают, но т.к. они ищутся локальными методами и вычислительные схемы различны, параллельное решение обеих задач представляет интерес как в плане использования решения, полученного в рамках одной модели для получения уточненного решения в рамках другой, так и в плане сравнения результатов.

Для любой вектор-функции $f \in L^4(X, \mu; R^n)$ обозначим

$$\phi_i^{(k)}(f) = \int_X f_i^k(w) d\mu(w) \quad 1 \leq i \leq n; k=1,2,3,4 \quad (5)$$

$$\pi_{ij}^{(k,1)}(f) = \int_X f_i^k(w) f_j^1(w) d\mu(w) \quad 1 \leq i, j \leq n, i \neq j; k=1,2,3; \quad (6) \\ 1=1,2$$

$$\pi_{ijl}^{(k)}(f) = \int_X f_i^k(w) f_j^1(w) f_l^1(w) d\mu(w) \quad 1 \leq i, j, l \leq n, i \neq j, i \neq l, \quad (7) \\ j \neq l; k=1,2$$

$$\varphi_{i_1 j_1 l s}(f) = \int_X f_{i_1}(w) f_{j_1}(w) f_l(w) f_s(w) d\mu(w) \quad \begin{matrix} 1 \leq i_1, j_1, l, s \leq n, & i_1 \neq j_1, \\ & i_1 \neq l, & i_1 \neq s, & j_1 \neq l, & j_1 \neq s, & l \neq s \end{matrix} \quad (8)$$

$$q_j(f; z) = \int_X r(z, w) f_j(w) d\mu(w) \quad 1 \leq j \leq n \quad (9)$$

Заметим, что $\varphi_{ij}^{(k,l)}(f) = \varphi_{ji}^{(l,k)}(f)$, $\varphi_{ijl}^{(1)}(f) = \varphi_{i_1 i_2 i_3}^{(1)}(f)$, $\varphi_{i_1 j_1 l s}(f) = \varphi_{j_1 j_2 j_3 j_4}(f)$, где (i_1, i_2, i_3) , (j_1, j_2, j_3, j_4) соответственно произвольные перестановки чисел (i, j, l) и (i, j, l, s) .

Положим

$$r(z) = \int_X r(z, w) d\mu(w).$$

В дальнейшем, как обычно, будем предполагать симметричность функции близости r . Это не умаляет общности результатов, однако упрощает их формулировку.

ТЕОРЕМА 1. 1. Если вектор-функция $t \in L^{(n)}$ является решением проекционной модели инвариантного шкалирования в смысле задачи $\Pi^{(n)}_1$, то для п.в. $z \in (X, \mu)$ имеет место система равенств

$$\begin{cases} K_j^{(n)}(t; z) = 0 & 1 \leq j \leq n \\ \sum_{i=1}^n \{ \sigma_i^{(4)}(t) + 3[\sigma_i^{(2)}(t)]^2 \} + 2 \sum_{j>1}^n [\sigma_i^{(2)}(t) \sigma_j^{(2)}(t) + \\ + 2[\varphi_{ij}^{(1,1)}(t)]^2 + \varphi_{ij}^{(2,2)}(t)] = \frac{\|r\|^2}{2}, \end{cases} \quad (10)$$

где

$$\begin{aligned} K_j^{(n)}(t; z) = & t_j^3(z)(r, \rho_t) + t_j(z) \{ [3\sigma_j^{(2)}(t) + \\ & + \sum_{i=1}^n (t_i^2(z) + \sigma_i^{(2)}(t))] (r, \rho_t) - \|r\|^2 \bar{r}(z) \} + \\ & + \{ [-\sigma_j^{(3)}(t) + \sum_{\substack{i=1 \\ i \neq j}}^n (2t_i(z) \varphi_{ij}^{(1,1)}(t) - \\ & - \varphi_{ij}^{(2,1)}(t))] (r, \rho_t) + \|r\|^2 q_j(t; z) \}. \end{aligned}$$

2. Если вектор-функция $t \in L^{(n)}$ является решением проекционной модели инвариантного шкалирования в смысле задачи $\Pi^{(n)}_2$, то для п.в. $z \in (X, \mu)$ имеет место система равенств

$$\begin{cases} \bar{K}_j^{(n)}(t; \lambda; z) = 0 & 1 \leq j \leq n \\ \sum_{i=1}^n \{ \sigma_i^{(4)}(t) + 3[\sigma_i^{(2)}(t)]^2 \} + 2 \sum_{j>1}^n [\sigma_i^{(2)}(t) \sigma_j^{(2)}(t) + \\ + 2[t_{ij}^{(1,1)}(t)]^2 + t_{ij}^{(2,2)}(t)] = \frac{\|r\|^2}{2}, \end{cases} \quad (11)$$

где

$$\begin{aligned} \bar{K}_j^{(n)}(t; \lambda; z) = & 2\lambda t_j^3(z) + t_j(z) \{ 2\lambda [3\sigma_j^{(2)}(t) + \\ & + \sum_{\substack{i=1 \\ i \neq j}}^n (t_i^2(z) + \sigma_i^{(2)}(t))] - \bar{r}(z) \} + \\ & + \{ 2\lambda [-\sigma_j^{(3)}(t) + \sum_{\substack{i=1 \\ i \neq j}}^n (2t_i(z) t_{ij}^{(1,1)}(t) - \\ & - t_{ij}^{(2,1)}(t))] + q_j(t; z) \}, \end{aligned}$$

причем параметр λ является множителем Лагранжа и при выполнении системы равенств (11) имеет место равенство

$$\lambda = \frac{\psi_{\Pi Y}^{(n)}(t)}{2\|r\|^2}. \quad (12)$$

Системы уравнений (10) и (11) будем называть системами уравнений многомерного проекционного инвариантного шкалирования соответственно в смысле модели $\Pi^{(n)}_1$ или $\Pi^{(n)}_2$.

Рассмотрим задачи $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$, когда функция $r(z, w)$ вырождена, т.е. допускает конечное представление

$$r(z, w) = \sum_{v=1}^q P_v(z) Q_v(w), \quad (13)$$

где $P_v, Q_v \in L^{\infty}(X, \mu)$, $1 \leq v \leq q$. Представление (13) имеет место для ряда важных функций близости, встречающихся в приложениях при анализе как количественной так и качественной информации (например, для взвешенной евклидовой функции близости [1]).

Положим

$$\alpha_v = \int_X P_v(w) d\mu(w) \quad 1 \leq v \leq q$$

$$\beta_v = \int_X Q_v(w) d\mu(w) \quad 1 \leq v \leq q$$

Для любой вектор-функции $f \in L^4(X, \mu; R^n)$ определим два набора параметров

$$s_f^{(1)} = \left\{ \left\{ \epsilon_i^{(k)}(f) \right\}_{\substack{1 \leq i \leq n \\ k=2,3,4}}, \left\{ \gamma_{ij}^{(k,l)}(f) \right\}_{\substack{1 \leq i, j \leq n; i \neq j \\ k=1,2, l=1,2}} \right\},$$

$$\left\{ \left\{ \theta_{iv}^{(k)}(f) \right\}_{\substack{1 \leq i \leq n \\ k=1,2}}, \left\{ \eta_{iv}(f) \right\}_{\substack{1 \leq i \leq n \\ 1 \leq v \leq q}} \right\},$$

$$s_f^{(2)} = \left\{ \left\{ \epsilon_i^{(k)}(f) \right\}_{\substack{1 \leq i \leq n \\ k=2,3,4}}, \left\{ \gamma_{ij}^{(k,l)}(f) \right\}_{\substack{1 \leq i, j \leq n; i \neq j \\ k=1,2, l=1,2}} \right\},$$

$$\left\{ \left\{ \theta_{iv}^{(1)}(f) \right\}_{\substack{1 \leq i \leq n \\ 1 \leq v \leq q}} \right\},$$

где $\epsilon_i^{(k)}(f)$, $\gamma_{ij}^{(k,l)}(f)$ определяются равенствами (5) - (6) и

$$\theta_{iv}^{(k)}(f) = \int_X f_i^{(k)}(w) Q_v(w) d\mu(w) \quad 1 \leq i \leq n; \quad 1 \leq v \leq q; \quad k=1,2 \quad (14)$$

$$\eta_{iv}(f) = \int_X f_i(w) P_v(w) d\mu(w) \quad 1 \leq i \leq n; \quad 1 \leq v \leq q \quad (15)$$

Рассмотрим две совокупности $s^{(1)}$ и $s^{(2)}$ наборов соответственно

$$s^{(1)} = \left\{ \left\{ \sigma_i^{(k)} \right\}_{\substack{1 \leq i \leq n \\ k=2,3,4}}, \left\{ \tau_{ij}^{(k,k)} \right\}_{\substack{1 \leq i < j \leq n \\ k=1,2}}, \left\{ \tau_{ij}^{(2,1)} \right\}_{\substack{1 \leq i, j \leq n \\ i \neq j}} \right\},$$

$$\left\{ \theta_{iv}^{(k)} \right\}_{\substack{1 \leq i \leq n, \\ k=1,2}}, \left\{ \eta_{iv} \right\}_{\substack{1 \leq i \leq n \\ 1 \leq v \leq q}},$$

$$s^{(2)} = \left\{ \left\{ \sigma_i^{(k)} \right\}_{\substack{1 \leq i \leq n \\ k=2,3,4}}, \left\{ \tau_{ij}^{(k,k)} \right\}_{\substack{1 \leq i < j \leq n \\ k=1,2}}, \left\{ \tau_{ij}^{(2,1)} \right\}_{\substack{1 \leq i, j \leq n \\ i \neq j}} \right\},$$

$$\left\{ \theta_{iv}^{(1)} \right\}_{\substack{1 \leq i \leq n \\ 1 \leq v \leq q}}$$

вещественных чисел $\sigma_i^{(k)}$, $\tau_{ij}^{(k,1)}$, $\theta_{iv}^{(k)}$, η_{iv} , удовлетворяющих условию

$$\sum_{i=1}^n [\sigma_i^{(4)} + 3(\sigma_i^{(2)})^2] + 2 \sum_{\substack{j=1 \\ j \neq i}}^n [\sigma_1^{(2)} \sigma_j^{(2)} + 2(\tau_{ij}^{(1,1)})^2 + \tau_{ij}^{(2,2)}] = \frac{\|r\|^2}{2}.$$

Отметим, что в дальнейшем под $\tau_{ij}^{(1,2)}$ понимаем $\tau_{ji}^{(2,1)}$ и, если $i > j$, то под $\tau_{ij}^{(k,k)}$ понимаем всегда $\tau_{ji}^{(k,k)}$.

Очевидно, что для любого $t \in L^{(n)}$, являющегося решением системы уравнений (10) (или (11)) $s_t^{(1)} \in S^{(1)}$ (соответственно $s_t^{(2)} \in S^{(2)}$).

Определим отображения

$$M(s^{(1)}) : L^4(X, \mu; R^n) \rightarrow L^1(X, \mu; R^n) \text{ для любого } s^{(1)} \in S^{(1)},$$

$$M_\lambda(s^{(2)}) : L^4(X, \mu; R^n) \rightarrow L^1(X, \mu; R^n) \text{ для любого } s^{(2)} \in S^{(2)}$$

следующим образом. Для любых $f \in L^4(X, \mu; R^n)$ и $z \in X$ положим

$$M(s^{(1)})(f; z) = \left\{ [M(s^{(1)}; f)]_j(z) \right\}_{j=1}^n$$

$$M_\lambda(s^{(2)})(f; z) = \left\{ [M_\lambda(s^{(2)}; f)]_j(z) \right\}_{j=1}^n$$

$$[M(s^{(1)}; f)_j](z) = f_j^3(z) + f_j(z) \left[\sum_{\substack{i=1 \\ i \neq j}}^n f_i^2(z) - \delta_j^{(0)}(s^{(1)}; z) \right] + \\ + 2 \sum_{\substack{i=1 \\ i \neq j}}^n f_i(z) \eta_{ij}^{(1,1)} + \delta_j(s^{(1)}; z) \quad (16)$$

$$[M_\lambda(s^{(2)}; f)_j](z) = f_j^3(z) + f_j(z) \left[\sum_{\substack{i=1 \\ i \neq j}}^n f_i^2(z) - \delta_{\lambda; j}^{(0)}(s^{(2)}; z) \right] + \\ + 2 \sum_{\substack{i=1 \\ i \neq j}}^n f_i(z) \eta_{ij}^{(1,1)} + \delta_{\lambda; j}(s^{(2)}; z), \quad (17)$$

где

$$\delta_j^{(0)}(s^{(1)}; z) = \frac{\|x\|^2}{\varepsilon} \sum_{v=1}^q \beta_v P_v(z) - \sum_{i=1}^n \sigma_i^{(2)} - 2\sigma_j^{(2)} \quad 1 \leq j \leq n \quad (18)$$

$$\delta_j(s^{(1)}; z) = \frac{\|x\|^2}{\varepsilon} \sum_{v=1}^q \theta_{jv}^{(1)} P_v(z) - \sum_{\substack{i=1 \\ i \neq j}}^n \eta_{ij}^{(2,1)} - \sigma_j^{(3)} \quad 1 \leq j \leq n \quad (19)$$

$$\varepsilon = 2 \sum_{v=1}^q \sum_{i=1}^n [\alpha_v \theta_{iv}^{(2)} - \theta_{iv}^{(1)} \eta_{iv}] \quad (20)$$

$$\delta_{\lambda; j}^{(0)}(s^{(2)}; z) = \frac{1}{2\lambda} \sum_{v=1}^q \beta_v P_v(z) - \sum_{i=1}^n \sigma_i^{(2)} - 2\sigma_j^{(2)} \quad 1 \leq j \leq n \quad (21)$$

$$\delta_{\lambda; j}(s^{(2)}; z) = \frac{1}{2\lambda} \sum_{v=1}^q \theta_{jv}^{(1)} P_v(z) - \sum_{\substack{i=1 \\ i \neq j}}^n \eta_{ij}^{(2,1)} - \sigma_j^{(3)} \quad 1 \leq j \leq n \quad (22)$$

Аналогично [2] для задач проекционного инвариантного шкалирования в смысле моделей $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$ вводится понятие области разрешимости, соответственно $S_*^{(1)}$ и при фиксированном значении параметра λ $S_*^{(2)}(\lambda)$. $S_*^{(1)}(S_*^{(2)}(\lambda))$ — это такое подмножество $S^{(1)}(S^{(2)})$, что для любого набора $s^{(1)} \in S_*^{(1)}$

$(s^{(2)} \in S_*^{(2)}(\lambda))$ существует единственная вектор-функция $t \in L^4(X, \mu; R^n)$, удовлетворяющая для п.в. $z \in (X, \mu)$ системе уравнений

$$[M(s^{(1)}; t)_j](z) = 0 \quad 1 \leq j \leq n \quad (23)$$

(соответственно

$$[M_{\lambda}(s^{(2)}; t)_j](z) = 0 \quad 1 \leq j \leq n). \quad (24)$$

В случае модели $\Pi^{(n)}_1$ в области разрешимости $S_*^{(1)}$ рассмотрим систему $n(2n+2+3q)$ уравнений относительно набора параметров $s^{(1)}$

$$\int_X [t_1(s^{(1)})](w) d\mu(w) = 0 \quad 1 \leq i \leq n \quad (25)$$

$$\int_X [t_1(s^{(1)})]^k(w) d\mu(w) = e_i^{(k)} \quad 1 \leq i \leq n; \quad k=2,3,4 \quad (26)$$

$$\int_X [t_1(s^{(1)})]^k(w) [t_j(s^{(1)})]^k(w) d\mu(w) = \eta_{ij}^{(k,k)} \quad \begin{matrix} 1 \leq i < j \leq n; \\ k=1,2 \end{matrix} \quad (27)$$

$$\int_X [t_1(s^{(1)})]^2(w) [t_j(s^{(1)})](w) d\mu(w) = \eta_{ij}^{(2,1)} \quad \begin{matrix} 1 \leq i, j \leq n, \\ i \neq j \end{matrix} \quad (28)$$

$$\int_X [t_1(s^{(1)})]^k(w) \varrho_v(w) d\mu(w) = \theta_{iv}^{(k)} \quad 1 \leq i \leq n; \quad 1 \leq v \leq q; \quad k=1,2 \quad (29)$$

$$\int_X [t_1(s^{(1)})](w) P_v(w) d\mu(w) = \eta_{1v} \quad 1 \leq i \leq n; \quad 1 \leq v \leq q, \quad (30)$$

где вектор-функция $[t(s^{(1)})](z)$ является решением системы n кубических уравнений (23), которая в области разрешимости имеет единственное вещественное решение относительно вектор-функции $t(z)$.

Так как суть задачи $\Pi^{(n)}_2$ состоит в максимизации функционала $\Psi^{(n)}_{\Pi}$, то в силу равенства (12) ищем наибольшее λ , при котором в области разрешимости $S^{(2)}_*(\lambda)$ система $n(2n+2+q)$ уравнений (25)–(28) и (29) при $k=1$ (везде в (25)–(29) делаем замену $s^{(1)}$ на $s^{(a)}$) относительно набора параметров $s^{(2)}$, где вектор-функция $[t(s^{(2)})](z)$ является единственным решением системы кубических уравнений (24) с параметром λ , имеет решение.

ТЕОРЕМА 2. 1. Если существует такая вектор-функция проекционного n -мерного инвариантного шкалирования в смысле модели $\Pi^{(n)}_1$ ($\Pi^{(n)}_2$) $t \in L^{(n)}$, что имеет место включение $s^{(1)}_t \in S^{(1)}_*$ ($s^{(2)}_t \in S^{(2)}_*(\lambda)$ при некотором λ), то система уравнений (25)–(30) ((25)–(28) и (29) при $k=1$ и этом λ) всегда имеет решение $s^{(1)}_* \in S^{(1)}_*$ ($s^{(2)}_* \in S^{(2)}_*(\lambda)$).

2. Для любого решения $s^{(1)}_* \in S^{(1)}_*$ (при фиксированном λ $s^{(2)}_* \in S^{(2)}_*(\lambda)$) системы (25)–(30) ((25)–(28) и (29) при $k=1$) вектор-функция $t(s^{(1)}_*)$ ($t(s^{(2)}_*)$), определенная системой уравнений (23) ((24)), является решением системы уравнений (10) ((11) при этом λ).

ЗАМЕЧАНИЕ 1. Так же как и в случае метрической модели любое решение системы (25)–(30), (23) ((25)–(28), (29) при $k=1$, (24)) относительно вектор-функции $t(z)$ порождает $2^n n!$ различных решений рассматриваемой проекционной задачи на подпространстве $L^{(n)}$, которые соответствуют замене знаков и различным перестановкам координатных функций вектор-функции $t(z)$.

Аналогично метрической модели, если $t(z)$ является решением системы уравнений (25)–(30), (23) ((25)–(28), (29) при

$k=1$, (24)), $A = (a_{1j})_{1,j=1}^n$ — ортогональная матрица, то в силу инвариантности проекционной модели относительно вращений $At(z)$ является решением соответствующей системы (25)–(30), (23) ((25)–(28), (29) при $k=1$, (24)), причем параметры новой системы $(\bar{\sigma}_1^{(k)}, \bar{\tau}_{1j}^{(k,1)}, \bar{\sigma}_{1v}^{(k)}, \bar{\eta}_{1v})$ (соответственно $\bar{\sigma}_1^{(k)}, \bar{\tau}_{1j}^{(k,1)}, \bar{\sigma}_{1v}^{(1)}$) пересчитываются по формулам

$$\bar{\sigma}_1^{(2)} = \sum_{j=1}^n a_{1j} \sigma_j^{(2)} + 2 \sum_{k>j=1}^n a_{1j} a_{1k} \tau_{jk}^{(1,1)} \quad 1 \leq i \leq n \quad (31)$$

$$\begin{aligned} \bar{\sigma}_1^{(3)} = & \sum_{j=1}^n a_{1j}^3 \sigma_j^{(3)} + 3 \sum_{j=1}^n \sum_{\substack{k=1 \\ k \neq j}}^n a_{1j}^2 a_{1k} \tau_{jk}^{(2,1)} + \\ & + 4 \sum_{1>k>j=1}^n a_{1j} a_{1k} a_{1l} \tau_{jkl}^{(1)}(t) \quad 1 \leq i \leq n \quad (32) \end{aligned}$$

$$\begin{aligned} \bar{\sigma}_1^{(4)} = & \sum_{j=1}^n a_{1j}^4 \sigma_j^{(4)} + 6 \sum_{k>j=1}^n a_{1j}^2 a_{1k}^2 \tau_{jk}^{(2,2)} + \\ & + 4 \sum_{j=1}^n \sum_{\substack{k=1 \\ k \neq j}}^n a_{1j}^3 a_{1k} \tau_{jk}^{(3,1)}(t) + 12 \sum_{j=1}^n \sum_{\substack{1>k=1 \\ k \neq j \\ l \neq j}}^n a_{1j}^2 a_{1k} a_{1l} \tau_{jkl}^{(2)}(t) + \\ & + 24 \sum_{p>1>k>j=1}^n a_{1j} a_{1k} a_{1l} a_{1p} \tau_{jklp}(t) \quad 1 \leq i \leq n \quad (33) \end{aligned}$$

$$\bar{\tau}_{1j}^{(1,1)} = \sum_{k=1}^n a_{1k} a_{jk} \sigma_k^{(2)} + \sum_{k=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n a_{1k} a_{jl} \tau_{kl}^{(1,1)} \quad 1 \leq i < j \leq n \quad (34)$$

$$\begin{aligned} \bar{\tau}_{1j}^{(2,2)} = & \sum_{k=1}^n a_{1k}^2 a_{jk}^2 \sigma_k^{(4)} + \sum_{k=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n (a_{1k}^2 a_{jl}^2 + 2 a_{1k} a_{1l} a_{jk} a_{jl}) \tau_{kl}^{(2,2)} + \\ & + 2 \sum_{z=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n (a_{1k}^2 a_{jk} a_{jl} + a_{1k} a_{1l} a_{jk}^2) \tau_{kl}^{(3,1)}(t) + \end{aligned}$$

$$\begin{aligned}
& + \sum_{k=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n \sum_{\substack{p=1 \\ p \neq k \\ p \neq l}}^n (a_{1k}^2 a_{jl} a_{jp} + a_{1l} a_{1p} a_{jk}^2 + 4a_{1k} a_{1l} a_{jk} a_{jp}) \xi_{klp}^{(2)}(t) + \\
& + 3 \sum_{k=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n \sum_{\substack{p=1 \\ p \neq k \\ p \neq l}}^n \sum_{\substack{m=1 \\ m \neq k \\ m \neq l \\ m \neq p}}^n a_{1k} a_{1l} a_{jp} a_{jm} \varphi_{klpm}(t) \quad 1 \leq i < j \leq n \quad (35)
\end{aligned}$$

$$\begin{aligned}
\bar{\eta}_{ij}^{(2,1)} &= \sum_{k=1}^n a_{1k}^2 a_{jk} \theta_k^{(3)} + \sum_{k=1}^n \sum_{\substack{l=1 \\ l \neq k}}^n (a_{1k}^2 a_{jl} + 2a_{1k} a_{1l} a_{jl}) \eta_{lk}^{(2,1)} + \\
& + 2 \sum_{k>l=1}^n \sum_{\substack{p=1 \\ p \neq k \\ p \neq l}}^n a_{1k} a_{1l} a_{jp} \xi_{klp}^{(1)}(t) \quad 1 \leq i, j \leq 1, \quad i \neq j \quad (36)
\end{aligned}$$

$$\bar{\theta}_{iv}^{(1)} = \sum_{j=1}^n a_{1j} \theta_{jv}^{(1)} \quad 1 \leq i \leq n; \quad 1 \leq v \leq q \quad (37)$$

$$\bar{\theta}_{iv}^{(2)} = \sum_{j=1}^n a_{1j}^2 \theta_{jv}^{(2)} + 2 \sum_{k>j=1}^n a_{1j} a_{1k} \psi_{jkv}(t) \quad 1 \leq i \leq n; \quad 1 \leq v \leq q \quad (38)$$

$$\bar{\eta}_{iv} = \sum_{j=1}^n a_{1j} \eta_{jv} \quad 1 \leq i \leq n; \quad 1 \leq v \leq q \quad (39)$$

где $\eta_{ij}^{(3,1)}(t)$, $\xi_{ijk}^{(k)}(t)$, $\varphi_{ijkl}(t)$ определяются формулами (6) - (8), а

$$\psi_{ijv}(t) = \int_X t_1(w) t_j(w) Q_v(w) d\mu(w) \quad 1 \leq i < j \leq n; \quad 1 \leq v \leq q.$$

Ниже приводятся достаточные условия, обеспечивающие принадлежность наборов параметров, порождаемых вектор-функцией проекционного шкалирования в рамках моделей $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$,

области разрешимости, и формулы, позволяющие вычислить значения вектор-функции шкалирования соответствующей этим наборам.

ТЕОРЕМА 3. Пусть $\tau(z) \in L^{(n)}$ является решением задачи $\Pi^{(n)}_1$ ($\Pi^{(n)}_2$). $A = (a_{ij})_{i,j=1}^n$ — какое-нибудь решение системы n^2 уравнений

$$\sum_{j=1}^n a_{ij}^2 = 1 \quad 1 \leq i \leq n \quad (40)$$

$$\sum_{k=1}^n a_{ik} a_{jk} = 0 \quad 1 \leq i < j \leq n \quad (41)$$

$$\sum_{k=1}^n a_{ik} [a_{jk} \sigma_k^{(2)}(\tau) + \sum_{l \neq k}^n a_{jl} t_{kl}^{(1,1)}(\tau)] = 0 \quad 1 \leq i < j \leq n \quad (42)$$

относительно вещественных a_{ij} .
(В случае модели $\Pi^{(n)}_2$ $\lambda = \frac{\Psi_{\Pi Y}^{(n)}(\tau)}{\|\tau\|^2}$.)

Тогда для того, чтобы набор параметров $s_t^{(1)}$ ($s_t^{(2)}$), порождаемый функцией проекционного n -мерного инвариантного шкалирования $t(z) = A\tau(z)$ принадлежал области разрешимости $s_*^{(1)}$ ($s_*^{(2)}(\lambda)$), достаточно, чтобы уравнение

$$T(z) \prod_{i=1}^n [T(z) - \Delta_j^{(0)}(z)]^2 - \sum_{i=1}^n \Delta_i^2(z) \prod_{\substack{j=1 \\ j \neq i}}^n [T(z) - \Delta_j^{(0)}(z)]^2 = 0, \quad (43)$$

где $\Delta_j^{(0)}(z) = \delta_j^{(0)}(s_t^{(1)}; z)$ ($\Delta_j^{(0)}(z) = \delta_{\lambda; j}^{(0)}(s_t^{(2)}; z)$),

$$\Delta_j(z) = \delta_j(s_t^{(1)}; z) \quad (\Delta_j(z) = \delta_{\lambda; j}(s_t^{(2)}; z)) \quad , \quad (44)$$

$\delta_j^{(0)}(s_t^{(1)}; z)$, $\delta_j(s_t^{(1)}; z)$, $\delta_{\lambda; j}^{(0)}(s_t^{(2)}; z)$, $\delta_{\lambda; j}(s_t^{(2)}; z)$

определяются равенствами (18) – (22) с помощью соотношений

(34)–(39) по вектор-функции $\varphi(z)$ и набору параметров $s_{\varphi}^{(1)}$ ($s_{\varphi}^{(2)}$), для п.в. $z \in (X, \mu)$ имело единственное вещественное решение, удовлетворяющее условию

$$\sum_{i=1}^n [\sigma_i^{(4)}(f) + 3(\sigma_i^{(2)}(f))^2] + 2 \sum_{j>1}^n [\sigma_i^{(2)}(f) \sigma_j^{(2)}(f) + 2(t_{ij}^{(1,1)}(f))^2 + t_{ij}^{(2,2)}(f)] = \frac{\|x\|^2}{2},$$

где $f(z) = (f_1(z), \dots, f_n(z))$,

$$f_1(z) = - \frac{\Delta_1(z)}{T(z) - \Delta_1^{(0)}(z)} \quad 1 \leq i \leq n,$$

а $T(z)$ – решение уравнения (43).

Более того, в этом случае для координатных функций вектор-функции $t(z) = (t_1(z), \dots, t_n(z))$ для п.в. $z \in (X, \mu)$ справедливы соотношения

$$t_1(z) = f_1(z) \quad 1 \leq i \leq n.$$

ЗАМЕЧАНИЕ 2. Условия (40), (41) в совокупности суть условие ортогональности матрицы A , а смысл условия (42) состоит в том, что поворот вектор-функции проекционного инвариантного шкалирования $\varphi(z)$, осуществляемый матрицей A гарантирует равенство нулю смешанных моментов вектор-функции $t(z) = A\varphi(z)$, т.е.

$$\int_X t_i(w) t_j(w) d\mu(w) = 0 \quad 1 \leq i, j \leq n, \quad i \neq j.$$

В случае статистической независимости строящихся шкал, т.е. в случае равенства нулю смешанных моментов вектор-функции проекционного инвариантного шкалирования $t(z)$, и выпол-

нения дополнительного требования нормировки на шкалы

$$\int_X t_i^2(w) d\mu(w) = \text{const} \quad 1 \leq i \leq n$$

справедлива

ТЕОРЕМА 4. Пусть решение задачи $\Pi^{(n)}_1$ ($\Pi^{(n)}_2$) таково, что

$$1) \quad t \in L^{(n)};$$

$$2) \quad t_{ij}^{(1,1)}(t) = 0, \quad 1 \leq i < j \leq n;$$

$$3) \quad \epsilon_1^{(2)}(t) = \text{const} = \epsilon_2(t), \quad 1 \leq i \leq n.$$

$$(В случае задачи \Pi^{(n)}_2 \quad \lambda = \frac{\psi^{(n)}(t)}{\|x\|^2}.)$$

Тогда для того, чтобы $s_{*}^{(1)} \in S_{*}^{(1)}$ ($s_t^{(2)} \in S_{*}^{(2)}(\lambda)$), достаточно, чтобы для п.в. $z \in (X, \mu)$

$$\Delta(z) < 3 \sqrt{\frac{\sum_{i=1}^n \Delta_i(z)}{4}}$$

$$\text{и} \quad \sum_{i=1}^n \epsilon_i^{(4)}(f) + n(2+n)[\epsilon_1^{(2)}(f)]^2 + 2 \sum_{j>1=1}^n t_{1j}^{(2,2)}(f) = \frac{\|x\|^2}{2},$$

где $\Delta_i(z)$ определяется соотношением (44),

$$\Delta(z) = \frac{\|x\|^2}{\epsilon} \sum_{v=1}^q \beta_v P_v(z) - (n+2)\epsilon_2(t),$$

$$(\Delta(z) = \frac{1}{2\lambda} \sum_{v=1}^q \beta_v P_v(z) - (n+2)\epsilon_2(t)),$$

$$f_1(z) = \frac{3 \Delta_1(z)}{5\Delta(z) + 3[(\frac{1}{2}v(z) + \sqrt{Q(z)})^{1/3} + (\frac{1}{2}v(z) - \sqrt{Q(z)})^{1/3}]}, \quad 1 \leq i \leq n,$$

$$Q(z) = \frac{2}{27} P^3(z) + \frac{1}{4} V^2(z),$$

$$P(z) = -\frac{1}{3} \Delta^2(z),$$

$$V(z) = \frac{2}{27} \Delta^3(z) - \sum_{i=1}^n \Delta_i^2(z),$$

причем в этом случае для координатных функций вектор-функции $t(z)$ для п.в. $z \in (X, \mu)$ справедливы соотношения

$$t_1(z) = f_1(z), \quad 1 \leq i \leq n.$$

Отметим, что в силу справедливости теоремы 4 так же как и в случае метрической модели [2] разумно в общем случае задачу проекционного многомерного инвариантного шкалирования в рамках обеих моделей $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$ рассматривать при дополнительном предположении выполнения условия статистической независимости строящихся шкал, т.е. для определения набора вектор-функции шкалирования, решать систему (25)–(30) ((25)–(28), (29) при $k=1$) при дополнительных условиях

$$f_{1j}^{(1,1)} = 0 \quad 1 \leq i < j \leq n \quad (45)$$

$$\sigma_1^{(2)} = \text{const} = \sigma_2 \quad 1 \leq i \leq n. \quad (46)$$

В этом случае условия 1), 2), 3) теоремы 4 оказываются выполненными. А тогда:

1) отпадает проблема решения системы (40)–(42), определяющей поворот начального приближения функции шкалирования, гарантирующий справедливость условия (45) для начального приближения решения системы (25)–(30) ((25)–(28), (29) при $k=1$);

2) упрощается проблема разрешимости, которая сводится к

условию существования единственного вещественного корня кубического уравнения, а не уравнения степени $2n + 1$ в общем случае;

3) вектор-функция шкалирования в рамках обеих моделей $\Pi^{(n)}_1$ и $\Pi^{(n)}_2$ выписывается явно через набор параметров, ее определяющий.

Доказательство всех теорем проводится аналогично доказательствам соответствующих теорем в [1-3].

Л и т е р а т у р а

1. Перекрест В.Т., Об одной модели одномерного шкалирования. Автоматика и телемеханика, № 2, 1980.
2. Хмельницкая А.Б., Об оптимизационных моделях многомерного инвариантного шкалирования. Сб. "Моделирование и применение ЭВМ в социологии" под ред. Ю.Н. Толстовой. ИСИ АН СССР, 1980.
3. Перекрест В.Т., Хмельницкая А.Б., Проекционные модели одномерного инвариантного шкалирования. Сб. "Труды конференции по применению математических методов в социологии" под ред. Э.П. Андреева. М., "Наука", 1980.
4. Хмельницкая А.Б., Об эквивалентности некоторых моделей инвариантного шкалирования. Сб. "Дискретный анализ и обработка эмпирической информации № 2" под ред. Б.Г. Миркина. СО АН СССР, Новосибирск, 1980.

С о д е р ж а н и е

Т. Мелс

Классификация и канонические модели сбалансированных полных факторных комплексов с произвольным числом факторов	3
---	---

Ю. Вардья, Т. Мелс

Новое семейство показателей статистической зависимости и его применения	24
---	----

А.-М. Парринг

Асимптотическое распределение коэффициентов полиномиальной регрессии	36
--	----

А. Левисто, Э. Тийт

Сравнение регрессии в разных подсовкупностях .	45
--	----

Р. Ээремаа

Классифицирование научных данных с получением частично покрываемых классов	55
--	----

Ю. Вилисмяз

Двухмерный статистический анализ данных	74
---	----

Л.М. Тоодинг

Эмпирическое исследование городской среды при помощи системы статистической обработки данных в ВЦ ТГУ	83
---	----

А.Б. Хмельницкая

О проекционных моделях многомерного инвариантного шкалирования	93
--	----

МЕТОДЫ МНОГОМЕРНОГО СТАТИСТИЧЕСКОГО АНАЛИЗА. Труды
вычислительного центра. Выпуск 44. На русском языке.
Тартуский государственный университет. ЭССР,
г. Тарту, ул. Оликооли, 18. Ответственный редактор
С. Коскед. Сдано в печать 11.12.79. Бумага печатная
30x42 1/4. Печ. листов 7,0 (условных 6,51). Учетно-
издат. листов 5,86. Тираж 500. МВ 09664. Типогра-
фия ТТУ, ЭССР, г. Тарту, ул. Пялсона, 14. Зак. №
1650. Цена 90 коп.